
Parsimony and Bayesian phylogenetics

Pablo A. Goloboff and Diego Pol

8.1 Introduction

Methods of phylogeny reconstruction are often divided into statistical methods (which require an explicit model of evolution) and non-statistical methods. Among methods with an explicit statistical justification, the most widely used are the methods of maximum likelihood, resulting from Felsenstein's (1973, 1981c) work, and more recently, Bayesian phylogenetic methods based on Monte Carlo Markov chains, following Li (1996), Mau and Newton (1997), and Larget and Simon (1999).

The aim of a statistically based method is to estimate tree topologies and values of possibly relevant parameters, as well as the uncertainty inherent in those estimations. A method that could do that with reasonable accuracy would be attractive indeed. It is often claimed that it is advantageous for a method to be based on a specific evolutionary model, because that allows incorporating into the analysis the 'knowledge' of the real world embodied in the model. Bayesian methods have become very prominent among model-based methods, in part because of computational advantages, and in part because they estimate the probability that a hypothesis is true, given the observations and model assumptions. Early work on phylogenetics suggested the desirability of probabilifying the falsehood or truth of hypotheses. This includes early papers by Farris (1973, 1977, 1978), who later reconsidered the question of whether phylogeny estimation is to be viewed as a statistical problem or not, and moved to the position that phylogenetic inference is best viewed in non-statistical terms (Farris 1983). When he first approached phylogeny as a statistical

problem, Farris (1973, p. 250) pointed out that the tree to be selected "should be the most probable tree on the basis of available data," and that (for tree T and data X) this probability (normally called *posterior probability*) can be calculated with Bayes' Theorem:

$$\Pr(T|X) = \frac{\Pr(X|T)\Pr(T)}{\Pr(X)}$$

where $\Pr(T)$ is the prior probability of the tree being analyzed (i.e. the probability, *a priori* of any observation, of the tree being the true one), the factor $\Pr(X|T)$ is the likelihood of the topology (i.e. the probability of the data, given the tree), and the denominator $\Pr(X)$ is the prior probability of the observed data (calculated as $\sum \Pr(X|T)\Pr(T)$ for all possible topologies). Farris (1973) noted that because the prior probability of each tree topology can be assumed to be the same (equal prior probabilities are usually called a *flat prior*), and because $\Pr(X)$ is fixed for the given observations, the choice, equivalent to parsimony, depends only on the likelihood of the tree. Farris (1973) developed a very general model, with minimal assumptions; under that model, the most likely tree is equivalent to the most-parsimonious tree. In the very same issue of *Systematic Zoology*, Felsenstein (1973) laid the basis for his subsequent developments of a very different model, with much more specific assumptions (including assumptions of Markovian evolution and Poisson substitution), conceived mostly as applicable to the evolution of DNA sequences. In the approach of Felsenstein (1981c) the values of parameters as well as branch lengths are jointly estimated

in order to maximize the likelihood function of a tree.

Bayesian approaches to phylogenetics have taken Felsenstein's methods a step further, and instead of producing point estimations of all parameters to maximize the likelihood, they have suggested integrating the likelihood across the different possible parameter values (i.e. branch lengths and substitution model parameters):

$$\Pr(X|T) = \int_{B_T} \int_{\Phi} \Pr(X|T, \beta_T, \varphi) f(\beta_T, \varphi) d\varphi d\beta_T$$

where B_T is the set of possible branch lengths (β_T) of topology T , Φ is the set of all possible substitution parameter values (φ) of the model, and $f(\beta_T, \varphi)$ is the prior distribution of these parameters. Both Farris (1973) and Felsenstein (1973) had considered such a type of integration desirable, but noted that a major problem with this approach is that it involves the calculus of a multidimensional integral for every possible topology, which is exceedingly complex and computationally demanding.

In order to overcome this problem, some researchers (e.g. Farris 1973; Hasegawa and Kishino 1989; Smouse and Li 1989) have attempted to compute the Bayesian posterior probability of a topology using the parameter values that maximize its likelihood factor (e.g. the maximum likelihood estimate of branch lengths). However, this approximation (as noted by Goloboff 2003, for the case of maximum likelihood) ignores an infinite number of additional hypotheses that result from alternative sets of branch lengths (or other parameter values) for that topology.

Others, instead, have suggested calculating the exact probabilities, integrating the likelihood of a topology across all possible sets of branch lengths (e.g. Rannala and Yang 1996) or other parameters (e.g. Sinsheimer *et al.* 1996). The complexity of this procedure, however, precludes its applicability to data sets of more than a few sequences, and therefore these methods were hardly ever used.

8.2 Markov chain Monte Carlo

Recently, three independent groups originally applied Markov chain Monte Carlo methods

(MCMC) to approximate the posterior probabilities of trees (Li 1996; Mau 1996; Mau and Newton 1997; Yang and Rannala 1997; Larget and Simon 1999; Mau *et al.* 1999; Newton *et al.* 1999; Li *et al.* 2000).

The idea in a MCMC is to make computationally feasible the integration of the posterior probabilities across the parameters of interest (e.g. topology, branch lengths, substitution parameters). The chain uses a proposal mechanism, which consists of gradual modifications from a starting point (ideally, randomly chosen), and it alternatively changes some parameter values (e.g. topology, branch lengths, substitution parameters), stochastically and aperiodically. These proposals or transitions are accepted with a probability given by the Metropolis–Hastings algorithm (Metropolis *et al.* 1953; Hastings 1970; see Larget and Simon 1999 or Huelsenbeck *et al.* 2002 for details) and the Markov chain proceeds until it reaches a stationary state.

If the Markov chain is irreducible (i.e. it is possible for the chain to visit every possible set of parameters and tree topologies) its stationary state converges to the joint posterior probability distribution of the parameters being modified (Tierney 1994). Thus, the frequency with which a given topology is visited in the Markov chain approximates its marginal posterior probability (Mau *et al.* 1999). Thus, the results of MCMC are directly interpreted in probabilistic terms; they can estimate the probability that a particular tree is the true tree for these sequences, conditional on the stochastic model of substitution (Li *et al.* 2000). Additionally, since the posterior probability distribution is simultaneously estimated, measures of uncertainty can be derived from the Markov chain.

The outcome of the stationary state of the MCMC is a set of phylogenetic trees (with their associated parameters). In phylogenetic applications of MCMC, the relative frequency of each topology (irrespective of branch length and substitution parameter values) is interpreted as its posterior probability (given the stochastic model and data). Therefore, it seems straightforward to take the topology with the highest posterior probability as the point estimate of the true topology. This was clearly recognized by several authors (Li 1996; Rannala and Yang 1996;

Yang and Rannala 1997; Larget and Simon 1999; Mau *et al.* 1999). The estimated tree was referred to as the maximum posterior probability (MAP) tree by Rannala and Yang (1996).

However, the availability of the approximation of the posterior distribution of trees also allows the evaluation of the variability of the estimates (e.g. topology or any other parameter integrated by the MCMC). As noted by Mau *et al.* (1999), summarizing the distribution of MCMC trees indeed presents a challenge and several methods have been proposed for such purpose.

For instance, Mau *et al.* (1999) and Rannala and Yang (1996) noted that a Bayesian “credible set” can be obtained as the collection of topologies having the sum of their posterior probabilities constrained to be no less than a specified value (e.g. 0.95). Li (1996) also considered alternative ways to estimate the posterior probability of the true phylogeny from the MCMC results, such as the use of the tree that has the minimum topological distance to the majority (e.g. 90%) of the MCMC trees, or the use of a majority-rule consensus of the set of topologies generated by MCMC.

In the latter option, the frequency of the clades has been interpreted by most Bayesian phylogeneticists (e.g. Huelsenbeck *et al.* 2002) as the posterior probability that the clade is true, following ideas of Newton *et al.* (1999) and Larget and Simon (1999). These authors propose summing the posterior probabilities of the trees in which each clade of the MAP is present as a way to summarize uncertainty in the tree topology estimate. This approach, which sums the frequency with which a particular clade appears in the Markov chain in order to estimate its posterior probability, is certainly the most commonly used way to summarize MCMC results. This approach is implemented in available software packages (e.g. MrBayes of Huelsenbeck and Ronquist 2001; BAMBE, of Simon and Larget 1998), and is frequently reported in empirical analyses using Bayesian methods. Here we will focus on some undesirable properties found on this frequently used option to summarize MCMC results. Other alternative ways to summarize the results are less frequently used, and they differ from this one in depending much more

on whether the chain has succeeded in finding the actual MAP(s). As the chain is not conceived as a search mechanism, but instead as a sampling mechanism, it is extremely unlikely that it will find the individual trees of maximum *a posteriori* probability, except in very small data sets.

8.3 Problems with estimations of monophyly by MCMC

In this section, the discussion will be within the realm of the rules and goals postulated by defenders of model-based methods. We also have other general concerns about model-based methods; these reflect a viewpoint not shared by Bayesians, and are therefore discussed in the following section. While the MCMC can be used to estimate any parameter of the evolutionary process, we are concerned here with the estimates that are relevant for phylogenetic studies: estimations of monophyly of groups. Other parameters, such as transition:transversion (ts:tv) ratios, while possibly the primary interest for other evolutionary studies, are only of secondary interest to the phylogeneticist. Part of the attraction of MCMC Bayesian methods is that the values estimated for those other parameters, such as ts:tv ratios, do not rely on estimation of a tree topology, an advantage for such studies which we do not dispute. However, the fact that our examples show that there are problems when the estimations of monophyly are carried out in a certain way suggests that establishing proper estimations from MCMC is far from automatic, and raises concerns about the validity of the inferences of those other parameters as well.

The most common approach to estimating probability of monophyly of a group *X* is by summing the posterior probabilities of all the trees where group *X* is monophyletic. This can be done for the groups present in the individual tree of highest posterior probability (as proposed in Larget and Simon 1999), or for each of the groups found in the analysis; these options make no difference for our argument.

Huelsenbeck *et al.* (2002, p. 674) claimed that, since “Bayesian inference is based on the likelihood function, it should inherit many of the nice statistical properties of the maximum-likelihood

method.” The “nice statistical property” for which likelihood has been held superior to parsimony is, quintessentially, statistical consistency. Statistical consistency has been proven for maximum likelihood, but only as a byproduct of the consistent estimation of the branch lengths between taxa (see Rogers 1997; Chang 1996; with discussion in Goloboff 2003). If the tree topologies are estimated without estimating branch lengths—integrating branch lengths for a given tree topology, as done in the Bayesian methods—then statistical consistency might be lost (as discussed in Goloboff 2003). And even if Bayesian analysis used optimal branch lengths (which would slow it down considerably), the fact that posterior probabilities of individual clades are estimated from sums of posterior probabilities of the trees having the clade still creates problems. So, the idea that Bayesian analysis should automatically “inherit the nice statistical properties of maximum likelihood” is no more than wishful thinking; Bayesian analysis with MCMC involves substantial modifications to maximum-likelihood.

Estimating the posterior probability for monophyly of a given group as the sum of posterior probabilities of the trees with that group may create serious problems, and it is easy to see why. Imagine that there is a single tree of highest likelihood¹, where group *X* is not present. Imagine that there are many trees of a likelihood only slightly inferior, where group *X* is monophyletic. The sum of the likelihoods of the trees with the group may exceed the likelihood of the one tree without the group, and then the method would conclude that the group has a relatively large probability of being monophyletic. While there are many situations under which such an asymmetry could occur, some of them are surprisingly simple. Consider the case of Fig. 8.1, a 25-taxon data set, with taxon *A* having only missing entries. The data determine a perfectly pectinate tree, except for the placement of *A*. The strict consensus for these data (analyzed with either parsimony or likelihood) is

an unresolved bush, which does express the fact that the monophyly of *no* group is actually supported by the data. Note that *A* can float in the skeleton tree of the remaining taxa; each of the 45 trees with alternative placements of *A* has exactly the same likelihood (and thus, under a flat prior on tree topologies, the same posterior probability). However, of all those trees, only two (*A* sister to *B*, or *A* sister to *C*) make the group *BC* non-monophyletic; the proportion of trees with the group *BC* monophyletic is thus $43/45 = 0.955$. That is almost exactly the posterior probability for monophyly of *BC* estimated by MrBayes (see Fig. 8.1; values on the branches are values reported by MrBayes, numbers above the branches are the frequencies of the groups in the 45 most-parsimonious trees). The group *BCD*, instead, is made non-monophyletic by two additional locations of taxon *A*, so it is monophyletic in a proportion of $41/45 = 0.911$. As one moves towards the middle of the tree, the proportions of locations which make the group non-monophyletic decreases: $6/45$ for group *BCDE*, $8/45$ for group *BCDEF*, etc. Past the middle of the tree, the proportions start increasing again. This is reflected almost exactly in the posterior probabilities reported by MrBayes. Since this is perfectly expected, the proposal mechanism used by MrBayes seems—at least for data sets as simple as this one—to provide a sample of the tree space adequate to estimate the sums of posterior probabilities for different groups; our criticism has nothing to do with sampling problems, but simply with the quantity that is being estimated. The (estimated) sum of posterior probabilities of the trees with and without a group provides a measure with no apparent utility. Using such a measure leads to the unfounded conclusion that, even when what we know about *A* is nothing, we can still estimate with some precision its placement in the tree! That location of *A* is determined rather by the priors on trees, but that means that the priors on groups are highly unequal. That an equal prior on trees may mean an unequal prior on groups has been discussed by Pickett and Randle (2005); Pickett and Randle also note that using flat priors on some aspects of a simulation may impose non-flat priors on other aspects. While the non-flat priors on groups (which undoubtedly exist) influence the posterior

¹ Whether the likelihood is calculated as the likelihood for optimal branch lengths, or the sum of the likelihoods for all the branch lengths for the given topology, makes no difference to our argument.

ROOT	AA
X	AA
W	GGGGAA
V	GGGGGGGAA
U	GGGGGGGGGAA
T	GGGGGGGGGGAAA
S	GGGGGGGGGGGAAA
R	GGGGGGGGGGGGAAA
Q	GGGGGGGGGGGGGAA
P	GGGGGGGGGGGGGGAAA
O	GGGGGGGGGGGGGGGAA
N	GGGGGGGGGGGGGGGGAAA
M	GGGGGGGGGGGGGGGGGAA
L	GGGGGGGGGGGGGGGGGGAAA
K	GGGGGGGGGGGGGGGGGGGAA
J	GGGGGGGGGGGGGGGGGGGGAAA
I	GGGGGGGGGGGGGGGGGGGGGAAA
H	GGGGGGGGGGGGGGGGGGGGGGAA
G	GGGGGGGGGGGGGGGGGGGGGGGAAA
F	GGGGGGGGGGGGGGGGGGGGGGGGAA
E	GGGGGGGGGGGGGGGGGGGGGGGGGAAA
D	GGGGGGGGGGGGGGGGGGGGGGGGGGAA
C	GGGGGGGGGGGGGGGGGGGGGGGGGGGAAA
B	GGGGGGGGGGGGGGGGGGGGGGGGGGGGAA
A	??

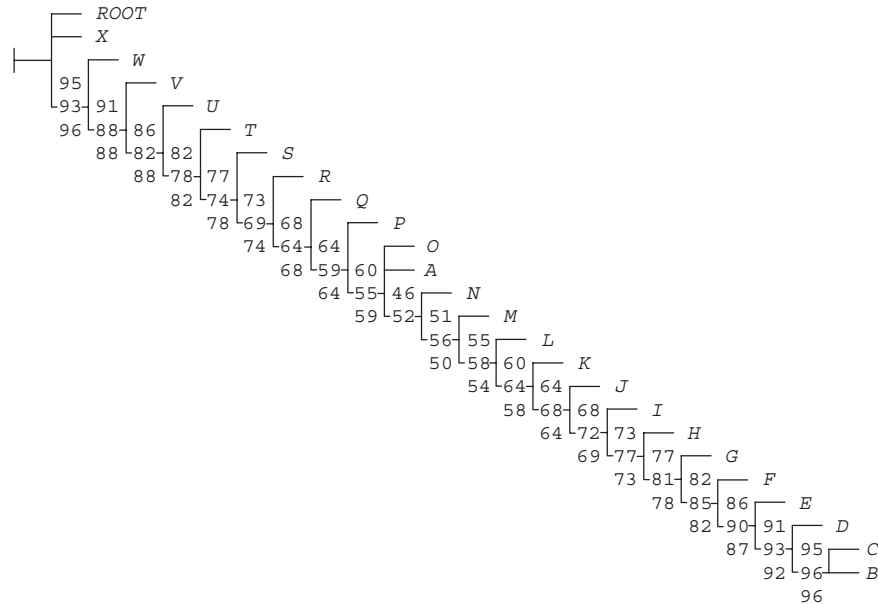


Figure 8.1 A data set with a taxon (A) scored only with missing entries. No group has any actual support, since the monophyly of any group can be violated at no cost. The numbers on the branches are the posterior probabilities of monophyly, estimated by MrBayes with 100 000 generations, using four chains, with a sampling frequency of 100, and a ‘burn-in’ of 250 (i.e. discarding the first 25 000 generations). The numbers above the branches indicate group frequency in the most-parsimonious (dichotomous) trees. The numbers below the branches show the bootstrap frequencies, as calculated by PAUP* (with 100 replications, analyzing each resampled data set with a branch-and-bound solution). Tree topology corresponds to the analysis with MrBayes.

probabilities reported by MrBayes, that is only part of the picture; the other aspect is the shape of the likelihood landscape, which is what our examples show.

Admittedly, the example of Fig. 8.1 is contrived in that no worker will attempt to analyze a matrix where a taxon is represented only by missing entries. But the same effects may come in much

more subtle flavors; for example, a sub-clade of a larger clade that can connect with different rootings (all with about the same likelihood) to the rest of a tree will produce, inside the clade with an undetermined root, the same effect observed for Fig. 8.1. This can also happen even for groups of a relatively large size (which non-flat priors on groups of different sizes do not easily explain), as in the example of Fig. 8.2, analyzed with MrBayes under the No Common Mechanism model (= parsimony). Under such a model, group *N–Z* is well supported by the data, and group *T–Z* is not: each of the characters that might support the monophyly of *T–Z* becomes an ambiguous synapomorphy when taxon *M* is the sister group of *N–Z* (so that there are trees of best fit that do not have *T–Z* as monophyletic). However, this happens only when *M* is the sister group of *N–Z*; for each of the other (numerous) possible locations of *M* in the rest of the tree, the group *T–Z* is required to provide the best fit to the data. The group *T–Z* is present in about 91% of the most-parsimonious

ROOT	AAAAAAAAA	AAAAA
A	AAAAAAAAA	AAAAA
B	AAAAAAAAA	AAAAA
C	AAAAAAAAA	AAAAA
D	AAAAAAAAA	AAAAA
E	AAAAAAAAA	AAAAA
F	AAAAAAAAA	AAAAA
G	AAAAAAAAA	AAAAA
H	AAAAAAAAA	AAAAA
I	AAAAAAAAA	AAAAA
J	AAAAAAAAA	AAAAA
K	AAAAAAAAA	AAAAA
L	AAAAAAAAA	AAAAA
M	AAAAAAAAA	GGGGG
N	GGGGGGGGG	AAAAA
O	GGGGGGGGG	AAAAA
P	GGGGGGGGG	AAAAA
Q	GGGGGGGGG	AAAAA
R	GGGGGGGGG	AAAAA
S	GGGGGGGGG	AAAAA
T	GGGGGGGGG	GGGGG
U	GGGGGGGGG	GGGGG
V	GGGGGGGGG	GGGGG
W	GGGGGGGGG	GGGGG
X	GGGGGGGGG	GGGGG
Y	GGGGGGGGG	GGGGG
Z	GGGGGGGGG	GGGGG

Figure 8.2 A data set with a group (*T–Z*) unsupported but found in many optimal trees, and thus with a high estimated posterior probability. See text for details.

trees for the data set². Not surprisingly, MrBayes reports unsupported group *T–Z* as strongly supported, with a posterior probability of 0.93.

The examples of Figs 8.1 and 8.2 were not derived from any model, and for this reason may perhaps be dismissed by Bayesians as being irrelevant. But the same effect can appear even in simulated data, where there are no violations of the model. The easiest way to produce the effect is to mimic the conditions of Fig. 8.1. For this, we used as the model tree a perfectly pectinate tree, as in Fig. 8.3, with taxa *A* and *B* forming a monophyletic group at the tip of the tree, and successive terminals appearing as successive sister groups. All the branches in the tree had a length of

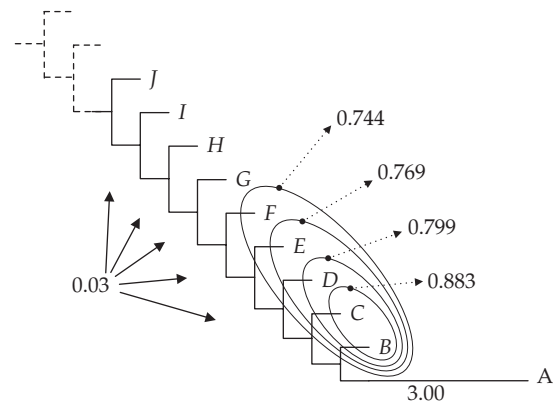


Figure 8.3 Tree shape used in the simulations (results reported in Fig. 8.4). Data were generated for trees with different numbers of taxa, using a Jukes–Cantor model. All branches were of length 0.03, except the branch leading to *A*, with a length of 3.0. The simulations generated 1 000 characters each. MrBayes analyses used 50 000 generations, with three chains, sampling every 50 generations, and a burn-in of 250 (i.e. discarding the first 12 500 generations). The posterior probabilities are shown for four incorrect groups (for 50 taxa, average posterior probability for 20 replications); note that the posterior probabilities decrease towards the middle of the tree, just as in Fig. 8.1.

² We calculated the frequency of group *T–Z* in most-parsimonious trees by taking a pseudo-random sample of 1000 most-parsimonious trees. We generated each by a Wagner tree where both the insertion and addition sequences were randomized (as implemented in TNT; Goloboff *et al.* 2004), followed by tree bisection and reconnection (TBR) branch swapping. Randomizing the insertion sequence means that, for each taxon to be added to form the Wagner tree, the pre-existing locations to insert the new taxon are tried in a random order; this eliminates bias in tree shapes in the resulting Wagner trees for poorly informative data.

0.03 (thus, a probability of no change along the branch of 0.978; we used a Jukes–Cantor model), except for the branch leading to taxon *A*, which was very long (with a length of 3; that is, a probability of no change of 0.287). The model tree was used to generate simulated data sets, with 1000 characters, for different numbers of taxa. Since *A* has a very long branch, it connects to the rest of the tree with about the same likelihood at every possible location. The effect is therefore the same as that of Fig. 8.1. In most of the simulations, MrBayes reports a high posterior probability that the group *BC* is monophyletic, which is in fact false. The estimated probability of monophyly of the wrong group *BC* actually increases with the number of taxa, since then the alternative locations of *A* that make a significant contribution to the sum of posterior probabilities for group *BC* also increases. Because there is significant variability in different simulated data sets, we used 20 replications for each of 5, 10, 20, 30, 40 and 50 taxa. The results are shown in Fig. 8.4. While there is of course some sampling error in our measurements, the trends evident in Fig. 8.4 make it clear that the high posterior probability attributed to the wrong

groups (often over 0.90) is not the effect of sampling error or lack of convergence in the chains. As the number of taxa increases, so does the apparent confidence on the false groups (the more so for the smaller groups), while the confidence on the true groups decreases (the more so for the smaller groups). Whereas the reference to smaller and larger groups makes sense in these examples, with pectinate trees, this does not mean that MCMC analysis will in general favor groups of a given size; the problem arises because of the relative differences in likelihood (= posterior probabilities, since we used a flat prior on trees) of those trees with and without each group, and this effect could potentially happen for groups of any size. These results could possibly derive as well from violations of the model, or from examining data for several genes where some of the taxa have not been sequenced for all the genes.

The examples are not intended to be realistic, but they show unequivocally that the estimations of posterior probabilities of individual groups may lead to grossly mistaken conclusions, and in real cases such an effect can easily be confounded by other factors. For the simulated examples, a taxon

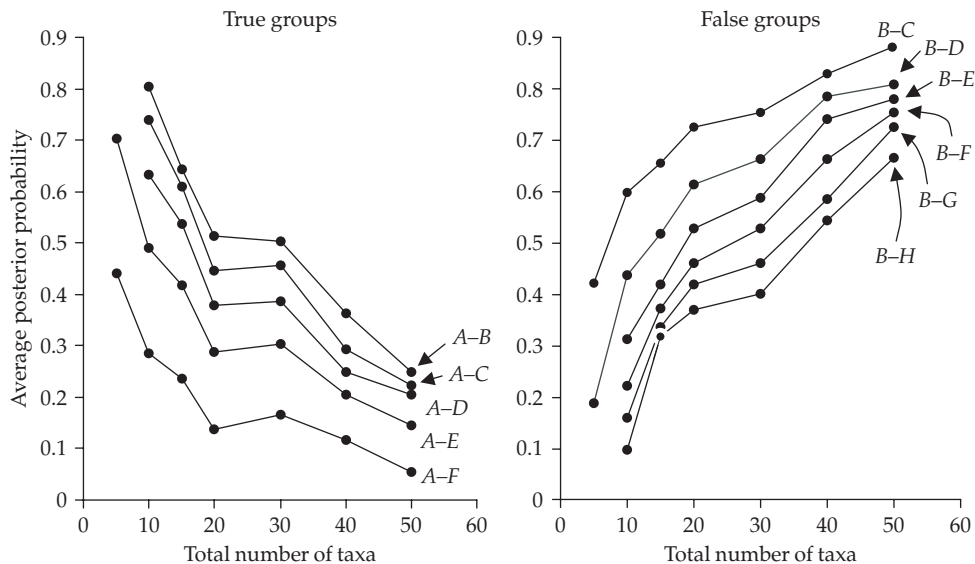


Figure 8.4 Results for the simulations, using the model tree shown in Fig. 8.3, for different numbers of taxa. As the number of taxa increases, so does the estimated posterior probability of the false groups (*BC*, *BCD*, etc.), the more so the smaller the group. All the averages reported correspond to 20 replications for each number of taxa.

with a branch as long as the branch leading to *A* cannot be confidently placed anywhere in the tree; every location will have roughly the same fit. The most serious problem faced by Bayesian analysis is not that it places *A* in some definite location (i.e. in the middle of the tree), but rather that it leads one to conclude that there is a very high probability that *A* is not placed as sister to *B*, which is the one true placement. A proper method should recognize, in cases like Fig. 8.3, that no conclusion is possible. Note that our criticism of Bayesian analysis, in this case, is not equivalent to Siddall's (1998) criticism of likelihood; Siddall (1998) criticized likelihood because in his simulations it separated long branches that were in fact sisters; Swofford *et al.* (2001) showed in their reply that, while likelihood indeed separates long sister branches (for small numbers of characters), the likelihoods of the alternative trees that place those long branches together is only slightly inferior, so that the maximum likelihood analysis actually implies that no decision is possible. That is not the case for the Bayesian results; they attribute a high probability to false groups that should at least be recognized as ambiguous.

8.4 Potential problems of the statistical approach

Statistically justified, model-based approaches to phylogeny have come to dominate the field in the last decade, but many authors still feel that those model-based justifications miss the mark. The controversy, not surprisingly, has often involved criticism and even misrepresentation from both sides. Among the topics on which the debate has centered are the questions of statistical consistency, the complexity of the evolutionary models used, the possible empirical basis of the evolutionary models on which the inferences are to be based, and whether epistemological considerations support the use of specific models of evolution.

The issue of consistency has been discussed mostly in relation to the likelihood vs. parsimony controversy (Steel *et al.* 1993; Siddall 1998; Farris 1999; Swofford *et al.* 2001). We consider that consistency, since it is relevant, at best, only under

unrealistic conditions (i.e. infinite, or at least massive amounts of data evolving under the same model³, with a perfect fit to the model), is not a very relevant property at the time of deciding among possible methods of phylogenetic inference. Even some statistically inclined phylogeneticists hold this point of view (e.g. Kim 1996; Sanderson and Kim 2000). The focus of this chapter is on Bayesian methods: Bayesians, with their claim that Bayesian analysis, being based on likelihood methods, should inherit the "nice statistical properties" of likelihood (see above), have adhered (implicitly, at least) to the notion that consistency is desirable. However, as we show later, the only feasible implementation of Bayesian phylogenetic analyses is likely to suffer from inconsistency.

The complexity of the inferential models used has also appeared in the likelihood vs. parsimony controversy. Although several likelihoodists (Goldman 1990; Steel and Penny 2000; Lewis 2001) had suggested that parsimony requires estimation of many more parameters than maximum likelihood, Goloboff (2003) reconsidered the problem and concluded that, if anything, parsimony requires estimation of fewer parameters than traditional maximum likelihood methods (a similar conclusion had been reached by Farris 1986, p. 22). Bayesian phylogenetic methods could in theory integrate uncertainty in some parameters during the MCMC (thus not requiring 'estimation' of those parameters). The problem is that the parameter space to be explored then becomes more complex, so that the chain would have to be run for much longer to insure convergence and an adequate sampling.

Finally, the epistemological questions about using evolutionary models and their empirical basis are perhaps the problems that have been less openly discussed in the literature. Several authors have presented the controversy in terms of which of the two approaches can be justified under Karl

³ Note that we say here "under the same model." While it is true that current-day techniques allow sampling of very long DNA sequences, the chances of all the sites still obeying to the same model decrease as the sequences become longer and include different genes or gene regions (as pointed out by Pol and Siddall 2001). The amount of data available for a given model will always be in the order of a few kilobases.

Popper's philosophy (Popper 1968). Most notable among recent philosophical defenses of model-based approaches is perhaps the paper by de Queiroz and Poe (2003). They argue that parsimony can be justified as a Popperian approach only by reference to specific models of evolution. de Queiroz and Poe (2003) say that it is not true that characters provide falsifiers of phylogenies (which Farris 1983 had used to characterize phylogenetic hypotheses as falsifiable), because a phylogeny cannot *per se* make impossible any particular character distribution. de Queiroz and Poe (2003) argue that, having disposed of other possible justifications, the only way to falsify a phylogeny is to show that it is less probable than its rival, and that parsimony can only be justified as Popperian if coupled with specific evolutionary models that specify those probabilities. But de Queiroz and Poe (2003) have not actually disposed of other possible justifications: they present only part of Farris' arguments. Farris (1983) had made it clear that no character could provide absolute falsification, and that the relationship between falsifier and hypothesis is purely logical. That is, if a given apparent character-state homology between two taxa is truly due to common ancestry, then it follows that a phylogeny that places them apart is truly false. Contra de Queiroz and Poe (2003), a strictly logical justification of parsimony, made without reference to a specific evolutionary model, is possible. Probabilistic models are necessary only to interpret the results of a parsimony analysis probabilistically; they are unnecessary otherwise.

The question of whether some apparent homologies are more probably truly homologous than others only enters the picture when we accept statistical justification and specific models. Contrary to some defenders of parsimony (e.g. Siddall and Kluge, 1997; Kluge 1997, 2001; Kluge, Chapter 2 of this volume), we do not argue that such a type of justification is philosophically and intrinsically flawed. Our concerns have to do with common sense, more than with philosophy. If one knew with certainty that sequence evolution is driven exclusively by a reduced set of parameters, and that those parameters remain very stable over time, then model-based methods of phylo-

geny reconstruction would be perfectly justified. Using those model-based methods would have the advantage that they make it possible to provide measures of uncertainty with a direct interpretation.

The alternative is considering that sequence evolution is driven by too many parameters, which may change too much over time, and that the samples (of sequences) we may expect to ever obtain are far below what could reasonably allow accurate inferences. Of course no one expects inferences that are 100% error-free; but the problem is how much is too much. Philosophical positions aside, many people who use parsimony place themselves at the 'too much' side of the scale, and tend to think that the probabilities estimated by using specific models are likely to be so far off that there is no point in trying to consider the results in terms of real probabilities. All we can expect is to simply provide the best explanation of the data, and it is best to remain silent about the probability of the resulting hypothesis being true. When de Queiroz and Poe (2003) claim that parsimony can be justified only by reference to some specific model, they mean "*parsimony as a statistical method can be justified only by reference to some specific model,*" which is true in itself, but then most proponents of parsimony do not view parsimony as attempting to provide the tree with the highest posterior probability: any attempt to provide a figure representing an actual probability, in the case of a process as complex as phylogeny, is no more reliable than a guess⁴. In this sense, the number of things that model-based methods try to estimate (statistically speaking) is much greater, and then it is natural that researchers with no previous experience in the field are attracted by estimation methods which are apparently omnipotent.

To some extent, the two aims—providing the best possible rationalization of the data by means of a phylogeny, or providing the best statistical

⁴ Measures of support such as the Bremer support (Bremer 1994; Goloboff and Farris 2001) or resampling (Farris *et al.* 1996; Goloboff *et al.* 2003b) are often interpreted as somehow measuring the truth content of the hypothesis, but this is not correct: all they measure is how much evidence supports the hypothesis.

estimation of the phylogeny—are both defensible in their own right. The difference is not only regarding whether a probabilistic model forms the basis for inferences or not; the difference is also about how the results are to be interpreted. Which aim a particular worker prefers and pursues may depend on many factors, actual personal interests, or even peer pressure, among others. But, fashions in science aside, the decision of whether using phylogenetic methods based either on models or pure logic depends also to a good extent on the dose of skepticism the researcher holds. Several defenders of model-based methods (e.g. Swofford *et al.* 2001) have suggested that, in different fields of science, the first approach is based on intuitive methods and that, as the field becomes mature, explicit statistical justifications replace the original intuitive ones. This claim is not strictly true (or testable, at least), and it is presented as if somehow the current use of statistically justified methods was evidence of maturity—when the alternative interpretation, namely that the use of statistical methods in phylogenetics is still premature and unjustified, may be much more reasonable to some workers. But how is one to decide whether the field is mature enough, or whether our knowledge of the mechanisms of evolution is detailed enough, to justify using those models? The answer to this need not be an all-or-none answer; there is instead a gray area between those who believe that our ignorance of evolutionary mechanisms is almost total (a view which some supporters of parsimony seem to hold), and those who believe that our knowledge is so complete as to guarantee even the most detailed inferences (as some likelihoodists and Bayesians seem to believe). At what particular point of this scale a particular worker finds himself/herself will depend on how he/she resolves a large number of subtle issues; such a decision requires reason and logic, but it cannot be accomplished with a statistical test. This, of course, is not to say that anything goes; for example, a model-based method that may produce incorrect estimations—even without violations of the model from which it is derived—is clearly to be avoided. Such is the case with estimations of posterior probabilities of monophyly by MCMC, as we have already seen.

Furthermore, whether the data can be modeled reasonably will depend on the nature of the data. While a Poisson model of substitution seems reasonable for some types of DNA sequence data, it seems unfounded to apply such a model to morphological data (although it has been attempted, by Lewis 2001). It is unfounded because there is very little ground to think that all characters of the given organisms have about the same probability of changing along a given branch of a tree, and that alternative states in morphological characters are like units turned on and off. Some vertebrates have mammary glands, and some arthropods have chelicerae, but within a given group of tetrapods, who could claim that the chances of gaining mammary glands are the same as—or even comparable to—the chances of gaining chelicerae? For genomic data, which include so many different types of transformations (insertions, deletions, translocations, inversions, etc.), postulating reasonable models is also very difficult or impossible. Therefore, for these types of data, only the parsimony approach has been used so far (starting from Sankoff and Blanchette 1998), for no other approach seems reasonable. Even the simplest case, insertions/deletions, presents a serious challenge to modeling; although some programs (like POY; Wheeler *et al.* 2003) have implemented “maximum likelihood models” for insertions/deletions; these Poisson “models” are based simply on attributing some probability to an insertion as a function of other parameters (like branch “length”). Wheeler *et al.* (personal communication⁵) present the likelihood methods in POY as “interpretive tools, without any necessary relationships to the actual process of change in nature.” They point out themselves that this has a meaning quite different from the use of Poisson models in DNA sequences: in those models, a base that is to replace another one exists outside the sequence and therefore, given a certain chance of replication error, there is a certain chance for each possible base to be inserted. Thus, Poisson

⁵ Wheeler, W., Aagesen, L., Arango, C., Faivovich, J., Grant, T., D’Haese, C., Janies, D., Smith, W., Varon, A. and Giribet, G. (unpublished manuscript). *Dynamic Homology and Phylogenetic Systematics: A Unified Approach using POY*.

substitution models are based on more than just attributing arbitrary probabilities of change to events; the probabilities postulated by those models are based on some knowledge of the mechanisms that govern the process of DNA substitution, at least in the absence of selection and constraints (whether the model is factually correct, of course, is a different matter, but it is plausible and coherent in itself). Gaps, on the other hand, are not units to be incorporated into a string of DNA being synthesized. A Poisson model for gaps, while it seems natural given the widespread use of Poisson models for DNA substitutions, may be totally inadequate. Other likelihood 'models' of insertions/deletions (e.g. Thorne *et al.* 1992; Miklós *et al.* 2004; see De Laet, Chapter 6 in this volume, for comments on these) do not use Poisson models, but still are based on arbitrarily assigning probabilities to the possible events. In those models, the final probability is no more real than are the figures obtained by parsimony analysis. On the other hand, analyzing sequences of unequal length by prealigning them and then discarding positions with gaps (a practice common among likelihoodists, and Bayesians) is probably even more inadequate, so that we are again in a situation where probabilities cannot really be assigned meaningfully. Much the same can be said of other types of chromosomal rearrangement.

8.5 Discussion

Strictly speaking, our simulations do not demonstrate that the estimations of posterior probabilities of individual groups produced by MrBayes are inconsistent. That would require either running data sets with infinite numbers of characters, or an analytical treatment of the multidimensional integral across all possible trees. Neither of those is possible. Admittedly, in cases like our simulations, as the number of characters increases, the difference in likelihood between the correct and alternative placements of the long branch increases. Eventually this difference might be so great as to make the likelihood of the individual best placement of the long branch (the correct one) higher than the sum of the likelihoods of the alternative (wrong) placements. However, as the number of taxa increases,

this situation becomes less and less likely, for the sum of likelihoods of the alternative placements increases as well. So, it is hard to predict what would happen for infinite numbers of characters in cases with very large numbers of taxa. However, even if there is the potential for Bayesian estimations of monophyly to provide correct topological estimations for infinite numbers of characters, there is still the problem that Bayesian analysis claims to do much more than simply producing consistent estimations: it also claims to measure the degree of support of the conclusions, in a statistical sense. Our examples show that it does not.

Several recent papers (Suzuki *et al.* 2002; Alfaro *et al.* 2003; Cummings *et al.* 2003) have compared bootstrap and clade credibility values. In terms of the problem discussed above, some of those comparisons could never have been very informative, despite large amounts of computational effort. For example, the study of Cummings *et al.* (2003) used over 15 years of CPU time, but examined only data sets with four taxa. The problems pointed out here with Bayesian analyses can arise only with larger numbers of taxa, so Cummings *et al.*'s effort could never have led to discovery of those problems. Moreover, even for larger numbers of taxa, the problem pointed out here could not have been discovered by comparing posterior probabilities with the bootstrap values produced by PAUP* (the program used in essentially all published comparisons between bootstrap and Bayesian credibilities; Swofford 2002). When bootstrapping or jackknifing, in the case of multiple trees for a resampled matrix, PAUP* weights each group found according to its frequency. This produces exactly the same results as summing posterior probabilities of monophyly: groups that are very frequent in optimal or quasi-optimal trees always appear as highly supported, regardless of their actual support. Fig. 8.1 shows, below the branches, the bootstrap values estimated by PAUP*; they are almost exactly the same as Bayesian estimates. The relative (although not universal) agreement between bootstrap and Bayesian estimations has been taken as mutual confirmation, but in fact MrBayes and PAUP*'s implementation of bootstrapping and jackknifing share similar biases. Alternative implementations of resampling

methods (such as the one in TNT; see Goloboff *et al.* 2004) avoid this problem by producing the strict consensus for each resampled matrix.

What happens with other ways to summarize the results of a MCMC? As noted above, they depend on whether the chain succeeded in finding the individual trees of highest posterior probability. For larger numbers of taxa, it is extremely unlikely that the chain will ever pass through the optimal tree(s), let alone pass through the optimal tree(s) enough times to estimate their posterior probability with any accuracy. Although tree bisection reconnection (TBR) forms the basis for both tree search and MCMC algorithms, rearrangements leading to worse trees are often accepted under MCMC, while they are normally rejected in a tree search. For equivalent numbers of rearrangements, then, a tree search (specially one combining different algorithms, like the methods in TNT; see Goloboff 1999, for details) is much more likely than a MCMC to find an optimal tree. Even if MCMC and a tree search had the same chances of finding an optimal tree for the same number of rearrangements, the numbers of rearrangements required to find optimal trees during a search cannot ever be achieved in Bayesian analyses. For example, in the case of a relatively small matrix of 84 taxa (from Goloboff 1995), TNT requires at least 5–10 million rearrangements to produce the first hit to minimum length. For Chase *et al.*'s (1993) data set (*zilla*, 500 taxa), it takes TNT an average of about 500 million rearrangements to find an optimal tree for the first time⁶; for the 854 taxa used in Goloboff (1999), it takes about 5000 million rearrangements (about 18 min in an 800 MHz

machine). Running 5 000 000 000 generations of a MCMC is impossible (in practical terms).

Suppose anyway that the chain succeeds in finding the trees of maximum *a posteriori* probability a certain number of times. The posterior probability of each individual tree will thus be negligible (the more so the more taxa are included in the analysis). In our view, such a low posterior probability is perfectly reasonable, and illustrates the fact that the statistical significance of phylogenetic conclusions cannot be meaningfully assessed in real cases. But statistically minded phylogeneticists will likely show continued interest in making probabilities more robust, i.e. in producing more 'acceptable' values. The alternative is to identify a credible set of trees. A strict consensus of the credible set of trees may contain exclusively well-supported groups, but only to the extent that the chain was run for long enough to find some trees that are relatively close to optimal trees. For the simulations carried out here (small numbers of taxa, very clean data without violations of the model, chains quickly converging), credibility sets of 90% still display, in many cases, false groups. Only running very large numbers of generations would avoid that problem, but in the case of larger data sets this will be impossible.

8.6 Acknowledgments

It is our pleasure to contribute to a book in the honor of James S. Farris, whose work in the field of phylogenetic systematics we greatly admire and respect. We also thank comments from Victor Albert, Jan De Laet, Mark Simmons, and Ward Wheeler. Financial support from FONCyT and CONICET (PAG), and The American Museum of Natural History (DP), is gratefully acknowledged.

⁶ Note that the figure of 500 million trees corresponds to analyses using sectorial searches, tree drifting, and tree fusing. Davis *et al.* (Chapter 7) report the numbers of rearrangements required to find optimal trees for *zilla* using only TBR; those numbers are much larger.

References

- Addario-Berry, L., Chor, B., Hallett, M., Lagergren, J., Panconesi, A. and Wareham, T. (2004). Ancestral maximum likelihood of phylogenetic trees is hard. *J. Bioinf. Comp. Biol.* **2**: 257–271.
- Aguinaldo, A.M., Turbeville, J.M., Linford, L.S., Rivera, M.C., Garey, J.R., Raff, R.A. and Lake, J.A. (1997). Evidence for a clade of nematodes, arthropods and other moulting animals. *Nature* **387**: 489–493.
- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *Second International Symposium on Information Theory*, Tsahkadsor, Armenia, USSR, September 2–8, 1971 (eds B.N. Petrov and F. Csaaki), pp. 267–281. Budapest, Akademiai Kiado.
- Albert, V.A. and Mishler, B.D. (1992). On the rationale and utility of weighting nucleotide sequence data. *Cladistics* **8**: 73–83.
- Albert, V.A., Mishler, B.D. and Chase, M.W. (1992). Character-state weighting for restriction site data in phylogenetic reconstruction, with an example from chloroplast DNA. In *Molecular Systematics of Plants*. (eds P.S. Soltis, D.E. Soltis, and J.J. Doyle), pp. 369–403. New York, Chapman and Hall.
- Albert, V.A., Chase, M.W. and Mishler, B.D. (1993). Character-state weighting for cladistic analysis of protein-coding DNA sequences. *Ann. Missouri Bot. Gard.* **80**: 752–766.
- Albert, V.A., Backlund, A., Bremer, K., Chase, M.W., Manhart, J.R., Mishler, B.D. and K.C. Nixon. (1994). Functional constraints and *rbcL* evidence for land plant phylogeny. *Ann. Missouri Bot. Gard.* **81**: 534–567.
- Alfaro, M., Zoller, S. and Lutzoni, F. (2003). Bayes or bootstrap? A simulation study comparing the performance of Bayesian Markov chain Monte Carlo sampling and bootstrapping in assessing phylogenetic confidence. *Mol. Biol. Evol.* **20**: 255–266.
- Altschul, S.F. (1989). Gap costs for multiple alignment. *J. Theor. Biol.* **138**: 297–309.
- Aravind, L., Watanabe, H., Lipman, D.J. and Koonin, E.V. (2000). Lineage-specific loss and divergence of functionally linked genes in eukaryotes. *Proc. Natl. Acad. Sci. USA* **97**: 11319–11324.
- Ariew, A. (1998). Are probabilities necessary for evolutionary explanations? *Biol. Philos.* **13**: 245–253.
- Arvestad, L., Berglund, A.C., Lagergren, J. and Sennblad, B. (2003). Bayesian gene/species tree reconciliation and orthology analysis using MCMC. *Bioinformatics* **19** Suppl. 1: i7–i15.
- Avise, J.C. (1989). Gene trees and organismal histories: a phylogenetic approach to population biology. *Evolution* **43**: 1192–1208.
- Avise, J.C. (2000). *Phylogeography: The History and Formation of Species*. Cambridge, MA, Harvard University Press.
- Bach, E. (1981). On time, tense and aspect: An essay in English metaphysics. In *Radical Pragmatics* (ed. P. Cole), pp. 63–81. New York, Academic Press.
- Baker, A. (2003). Quantitative parsimony and explanatory power. *Br. J. Phil. Sci.* **54**: 245–259.
- Bandelt, H.J., Forster, P., Sykes, B.C. and Richards, M.B. (1995). Mitochondrial portraits of human populations using median networks. *Genetics* **141**: 743–753.
- Barnes, E.C. (2000). Ockham's Razor and the anti-superfluity principle. *Erkenntnis* **53**: 353–374.
- Barrett, M., Donoghue M.J. and Sober, E. (1991). Against consensus. *Syst. Zool.* **40**: 486–493.
- Barry, D. and Hartigan, J.A. (1987). Statistical analysis of hominoid molecular evolution. *Stat. Sci.* **2**: 191–210.
- Benner, S.A., Trabesinger, N. and Scheiber, D. (1998). Post-genomic science: Converting primary sequence into physiological function. *Adv. Enzyme Regul.* **38**: 155–190.
- Benner, S.A., Chamberlin, S.G., Liberles, D.A., Govindarajan, S. and Knecht, L. (2000). Functional inferences from reconstructed evolutionary biology involving rectified databases — an evolutionarily grounded approach to functional genomics. *Res. Microbiol.* **151**: 97–106.
- Benson, D.A., Karsch-Mizrachi, I., Lipman, D.J., Ostell, J. and Wheeler, D.L. (2004). GenBank: update. *Nucleic Acids Res.* **32**: D23–D26.
- Bi, S., Garilova, O., Gong, D.W., Mason, M.M. and Reitman, M. (1997). Identification of a placental enhancer for the human leptin gene. *J. Biol. Chem.* **272**: 30583–30588.

- Blackburn, D.G. (1984). From whale toes to snake eyes: Comments on the reversibility of evolution. *Syst. Zool.* **33**: 241–245.
- Blair, J.E., Ikeo, K., Gojobori, T. and Hedges, S.B. (2002). The evolutionary position of nematodes. *BMC Evol. Biol.* **2**: 7.
- Bock, W.J. (1973). Philosophical foundations of classical evolutionary classification. *Syst. Zool.* **22**: 375–392.
- Boudet, N., Aubourg, S., Toffano-Nioche, C., Kreis, M. and Lecharny, A. (2001). Evolution of intron/exon structure of DEAD helicase family genes in *Arabidopsis*, *Caenorhabditis*, and *Drosophila*. *Genome Res.* **11**: 2101–2114.
- Boyd, R. (1991). Confirmation, semantics, and the interpretation of scientific theories. In *The Philosophy of Science* (eds R. Boyd, P. Gasper and J.D. Trout), pp. 3–35. Cambridge, MA, MIT Press.
- Brady, R.H. (1985). On the independence of systematics. *Cladistics* **1**: 113–126.
- Brady, R.H. (1994). Explanation, description, and the meaning of transformation in taxonomic evidence. In *Models in Phylogeny Reconstruction* (eds R.W. Scotland, D.J. Siebert and D.M. Williams), pp. 11–29 special vol. 52, Syst. Assoc. Oxford, Clarendon Press.
- Bremer, K. (1988). The limits of amino acid sequence data in angiosperm phylogenetic reconstruction. *Evolution* **42**: 795–803.
- Bremer, K. (1994). Branch support and tree stability. *Cladistics* **10**: 295–304.
- Brooks, D.R. (1981). Hennig's parasitological method: A proposed solution. *Syst. Zool.* **30**: 229–249.
- Brooks, D.R. (1988). Scaling effects in historical biogeography: A new view of space, time, and form. *Syst. Zool.* **37**: 237–244.
- Brower, A.V.Z. (2000). Homology and the inference of systematic relationships: Some historical and philosophical perspectives. In *Homology and Systematics: Coding Characters for Phylogenetic Analysis* (eds R. Scotland and R.T. Pennington), pp. 10–21. New York, Taylor and Francis.
- Brudno, M., Chapman, M., Götting, B., Batzoglou, S. and Morgenstern, B. (2003). Fast and sensitive multiple alignment of large genomic sequences. *BMC Bioinformatics* **4**: 66.
- Brudno, M., Poliakov, A., Salamov, A., Cooper, G.M., Sidow, A., Rubin, E.M., Solovyev, V., Batzoglou, S. and Dubchak, I. (2004). Automated whole-genome multiple alignment of rat, mouse, and human. *Genome Res.* **14**: 685–692.
- Bryant, D. (2000). A lower bound for the breakpoint phylogeny problem. In *Proceedings of the 11th Annual Symposium on Combinatorial Pattern Matching* (eds R. Giancarlo and D. Sankoff), pp. 235–247. London, Springer Verlag.
- Bryant, D. (2003). A classification of consensus methods for phylogenetics. In *Bioconsensus* (eds Janowitz, M.F., Lapointe, F.J., McMorris, F.R., Mirkin, B. and Roberts, F.S.), pp. 163–184. DIMACS Series in Discrete Mathematics and Theoretical Computer Science, vol. 61, Providence, RI, American Mathematical Society.
- Buneman, P. (1971). The recovery of trees from measures of dissimilarity. In *Mathematics in the Archaeological and Historical Sciences* (eds Hodson, F.R., Kendall, D.G. and Tautu, P.), pp. 387–395. Edinburgh, Edinburgh University Press.
- Burnham, K.P. and Anderson, D.R. (1998). *Model Selection and Inference. A Practical Information-theoretic Approach*. New York, Springer.
- Cameron, H.D. (1987). The upside-down cladogram: Problems in manuscript affiliation. In *Biological Metaphor and Cladistic Classification: An Interdisciplinary Perspective* (eds H.M. Hoenigswald and L.F. Wiener), pp. 227–242. Philadelphia, PA, University of Pennsylvania Press.
- Camin, J.H. and Sokal, R.R. (1965). A method for deducing branching sequences in phylogeny. *Evolution* **19**: 311–326.
- Carillo, H. and Lipman, D. (1988). The multiple sequence alignment problem in biology. *SIAM J. Appl. Math.* **48**: 1073–1082.
- Caroll, L. (1872). *Through the Looking Glass*. London, Macmillan.
- Carpenter, J.M. (2003). On “Molecular phylogeny of Vespidae (Hymenoptera) and the evolution of sociality in wasps”. *Am. Museum Novitates* **3389**: 1–20.
- Carter, M., Hendy, M.D., Penny, D., Székely, L.A. and Wormald, N.C. (1990). On the distribution of lengths of evolutionary trees. *SIAM J. Discrete Math.* **3**: 38–47.
- Cartmill, M. (1981). Hypothesis testing and phylogenetic reconstruction. *Z. Zool. Syst. Evol.-forsch.* **19**: 73–96.
- Chang, B.S.W., Jönsson K., Kazmi, M.A., Donoghue, M.J. and Sakmar T.P. (2002). Recreating a functional ancestral archosaur visual pigment. *Mol. Biol. Evol.* **19**: 1483–1489.
- Chang, J. (1996). Full reconstruction of Markov models on evolutionary trees: Identifiability and consistency. *Math. Biosc.* **137**: 51–73.
- Chang, J.T. and Kim, J. (1996). The measurement of homoplasy: A stochastic view. In *Homoplasy: the Recurrence of Similarity in Evolution* (eds M.J. Sanderson and L. Hufford), pp. 189–303. Academic Press.
- Chase, M.W., Soltis, D.E., Olmstead, R.G., Morgan, D., Les, D.H., Mishler, B.D., Duvall, M.R., Price, R.A., Hills, H.G., Qiu, Y.-L. et al. (1993). Phylogenetics of seed plants: An analysis of nucleotide sequences from the plastid gene *rbcL*. *Ann. Missouri Bot. Gard.* **40**: 528–580.

- Crichton, M. (1990). *Jurassic Park*. New York, Ballantine Books.
- Crick, F. (1968). The origin of the genetic code. *J. Mol. Biol.* **38**: 367–379.
- Crisci, J.V. and Stuessy, T.F. (1980). Determining primitive character states for phylogenetic reconstruction. *Syst. Bot.* **5**: 112–135.
- Cummings, M., Handley, S., Myers, D., Reed, D., Rokas, A. and Winka, K. (2003). Comparing bootstrap and posterior probability values in the four-taxon case. *Syst. Biol.* **52**: 477–487.
- Dacks, J.B. and Doolittle, W.F. (2001). Reconstructing/deconstructing the earliest eukaryotes: How comparative genomics can help. *Cell* **107**: 419–425.
- Dandekar, T., Snel, B., Huynen, M. and Bork, P. (1998). Conservation of gene order: A fingerprint of proteins that physically interact. *Trends Biochem. Sci.* **23**: 324–328.
- Darwin, C. (1859). *The Origin of Species by Means of Natural Selection, or the Preservation of Favoured Races in the Struggle for Life*. London, John Murray [1964. Facsimile of 1st edition. Cambridge, MA, Harvard University Press]
- Davids, W., Gamielien, J., Liberles, D.A. and Hide, W. (2002). Positive selection scanning reveals decoupling of enzymatic activities of carbamoyl phosphate synthetase in *H. pylori*. *J. Mol. Evol.* **54**: 458–464.
- Davidson, D. (1991). On the individuation of events. *Synthese* **86**: 229–254.
- Davis, J.I. and Nixon, K.C. (1992). Populations, genetic variation, and the delimitation of phylogenetic species. *Syst. Biol.* **41**: 421–435.
- Davis, J.I., Stevenson, D.W., Petersen, G., Seberg, O., Campbell, L.M., Freudenstein, J.V., Goldman, D.H., Hardy, C.R., Michelangeli, F.A., Simmons, M.P. et al. (2004). A phylogeny of the monocots, as inferred from *rbcL* and *atpA* sequence variation, and a comparison of methods for calculating jackknife and bootstrap values. *Syst. Bot.* **29**: 467–510.
- Dayhoff, M.O. and Eck, R.V. (1968). *Atlas of Protein Sequence and Structure*. 1967–68. Silver Spring, MD, National Biomed. Res. Foundation.
- Dayhoff, M.O. and Park, C.M. (1969). Cytochrome C: Building a phylogenetic hypothesis. In *Atlas of Protein Sequence and Structure*. 1969 (ed. M.O. Dayhoff), pp. 7–16 vol. 4. Silver Spring, MD, National. Biomed. Res. Foundation.
- DeBry, R.W. and Slade, N.A. (1985). Cladistic analysis of restriction endonuclease cleavage maps within a maximum-likelihood framework. *Syst. Zool.* **34**: 21–34.
- De Laet, J. (1997). *A Reconsideration of Three-Item Analysis, the Use of Implied Weights in Cladistics, and a Practical Application in Gentianaceae*. PhD Thesis, University of Leuven, Belgium.
- De Laet, J. (2003). Parsimony algorithms for characters that are inapplicable in some terminals (Abstract, 21st annual meeting of the Willi Hennig Society, Helsinki 2002). *Cladistics* **19**: 151.
- De Laet, J. (2004). When one and one is not two: Parsimony analysis of sequence data (Abstract, 22nd annual meeting of the Willi Hennig Society, New York 2003). *Cladistics* **20**: 81.
- De Laet, J. and Smets, E. (1998). On the three-taxon approach to parsimony analysis. *Cladistics* **14**: 363–381.
- De Laet, J. and Wheeler, W. (2003). *POY version 3.0.11* (Wheeler, Gladstein and De Laet, May 6 2003). Command line documentation. Available at ftp://ftp.amnh.org/pub/molecular/poy.
- de Pinna, M. C. C. (1991). Concepts and tests of homology in the cladistic paradigm. *Cladistics* **7**: 367–394.
- de Queiroz, K. (1992). Phylogenetic definitions and taxonomic philosophy. *Biol. Philos.* **7**: 295–313.
- de Queiroz, K. (1996). Including the characters of interest during tree reconstruction and the problems of circularity and bias in studies of character evolution. *Am. Nat.* **148**: 700–708.
- de Queiroz, K. and Poe, S. (2001). Philosophy and phylogenetic inference: A comparison of likelihood and parsimony methods in the context of Karl Popper's writings on corroboration. *Syst. Biol.* **50**: 305–321.
- de Queiroz, K. and Poe, S. (2003). Failed refutations: Further comments on parsimony and likelihood methods and their relationship to Popper's degree of corroboration. *Syst. Biol.* **52**: 352–367.
- Dezulian, T. and Steel, M. (2004). Phylogenetic closure operations and homoplasy-free evolution. In *Classification, Clustering, and Data Mining Applications* (Proceedings of the meeting of the International Federation of Classification Societies (IFCS) 2004). (eds D. Banks, L. House, F.R. McMorris, P. Arabie and W. Gaul), pp. 395–416. Springer-Verlag, Berlin.
- Dibb, N.J. and Newman, A.J. (1989). Evidence that introns arose at proto-splice sites. *EMBO J.* **8**: 2015–2021.
- Dollo, L. (1893). Le lois de l'évolution. *Bull. Soc. Belge Geol. Paleontol. d'Hydrolog.* **7**: 164–167.
- Donoghue, M.J. and Doyle, J.A. (1989). Phylogenetic analysis of angiosperms and the relationships of Hamamelidae. In *Evolution, Systematics and Fossil History of the Hamamelidae*, vol.1 (eds P. Crane and S. Blackmore), pp. 17–45. Oxford, Clarendon Press.
- Donoghue, M.J. and Sanderson, M.J. (1992). The suitability of molecular and morphological evidence in reconstructing plant phylogeny. In *Molecular Systematics of Plants*. (eds P.S. Soltis, D.E. Soltis and J.J. Doyle), pp. 340–368, New York, Chapman & Hall.

- Donoghue, M.J., Doyle, J.A., Gauthier, J.A., Kluge, A.G. and Rowe, T. (1989). The importance of fossils in phylogeny reconstruction. *Annu. Rev. Ecol. Syst.* **20**: 431–460.
- Doolittle, W.F. (1999). Lateral genomics. *Trends Cell. Biol.* **9**: M5–M8.
- Doolittle, W.F., Boucher, Y., Nesbo, C.L., Douady, C.J., Andersson, J.O. and Roger, A.J. (2003). How big is the iceberg of which organellar genes in nuclear genomes are but the tip? *Phil. Trans. R. Soc. Lond. B Biol. Sci.* **358**: 39–57.
- Duret, L., Mouchiroud, D. and Gouy, M. (1994). HOVERGEN: A database of homologous vertebrate genes. *Nucleic Acids Res.* **22**: 2360–2365.
- Edwards, A.W.F. (1972). *Likelihood*. Cambridge, Cambridge University Press.
- Edwards, A.W.F. and Cavalli-Sforza, L.L. (1963). The reconstruction of evolution. *Heredity* **18**: 553.
- Edwards, A.W.F. and Cavalli-Sforza, L.L. (1964). Reconstruction of evolutionary trees. In *Phenetic and Phylogenetic Classification* (eds V.H. Heywood and J. McNeill), pp. 67–76 no. 6. London, Syst. Assoc. Publ.
- Eisen, M.B., Spellman, P.T., Brown, P.O. and Botstein, D. (1998). Cluster analysis and display of genome-wide expression patterns. *Proc. Natl. Acad. Sci. USA* **95**: 14863–14868.
- Enard, W., Khaitovich, P., Klose, J., Zöllner, S., Heissig, F., Giavalisco, P., Nieselt-Struwe, K., Muchmore, E., Varki, A., Ravid, R. et al. (2002). Intra and interspecific variation in primate gene expression patterns. *Science* **296**: 340–343.
- Endo, T., Ikee, K. and Gojobori, T. (1996). Large-scale search for genes on which positive selection may operate. *Mol. Biol. Evol.* **13**: 685–690.
- Endress, P.K. (1994). *Diversity and Evolutionary Biology of Tropical Flowers*. Cambridge, Cambridge University Press.
- Erdős, P.L. and Székely, L.A. (1992). Evolutionary trees: An integer multicommodity maxflow–min-cut theorem. *Adv. Appl. Math.* **13**: 375–389.
- Erdős, P.L. and Székely, L.A. (1993). Counting bichromatic evolutionary trees. *Discrete Appl. Math.* **47**: 1–8.
- Erdős, P.L., Steel, M.A., Székely, L.A. and Warnow, T. (1999). A few logs suffice to build (almost) all trees (Part 1). *Random Struct Algorithms* **14**: 153–184.
- Excoffier, L. and Smouse, P.E. (1994). Using allele frequencies and geographic subdivision to reconstruct gene trees within a species: Molecular variance parsimony. *Genetics* **136**: 343–359.
- Farris, J.S. (1966). Estimation of conservatism of characters by constancy within biological populations. *Evolution* **20**: 319–334.
- Farris, J.S. (1967). The meaning of relationship and taxonomic procedure. *Syst. Zool.* **16**: 44–51.
- Farris, J.S. (1969). A successive approximations approach to character weighting. *Syst. Zool.* **18**: 374–385.
- Farris, J.S. (1970). Methods for computing Wagner trees. *Syst. Zool.* **19**: 83–92.
- Farris, J.S. (1972). Estimating phylogenetic trees from distance matrices. *Am. Nat.* **106**: 645–668.
- Farris, J.S. (1973a). On the use of the parsimony criterion for inferring phylogenetic trees. *Syst. Zool.* **22**: 250–256.
- Farris, J.S. (1973b). A probability model for inferring evolutionary trees. *Syst. Zool.* **22**: 250–256.
- Farris, J.S. (1977a). Phylogenetic analysis under Dollo's law. *Syst. Zool.* **26**: 77–88.
- Farris, J.S. (1977b). Some further comments on Le Quesne's methods. *Syst. Zool.* **26**: 220–223.
- Farris, J.S. (1978a). Inferring phylogenetic trees from chromosome inversion data. *Syst. Zool.* **27**: 275–284.
- Farris, J.S. (1978b). *Wagner78*. Published by the author.
- Farris, J.S. (1979). The information content of the phylogenetic system. *Syst. Zool.* **28**: 483–519.
- Farris, J.S. (1982a). Outgroups and parsimony. *Syst. Zool.* **31**: 328–334.
- Farris, J.S. (1982b). Simplicity and informativeness in systematics and phylogeny. *Syst. Zool.* **31**: 413–444.
- Farris, J.S. (1983). The logical basis of phylogenetic analysis. In *Advances in Cladistics*, volume 2: *Proceedings of the Second Meeting of the Willi Hennig Society*. (eds N. Platnick and V. Funk) pp. 7–36. New York, Columbia University Press.
- Farris, J.S. (1983/1994). The logical basis of phylogenetic analysis. In *Advances in Cladistics — Proceedings of the 2nd Annual Meeting of the Willi Hennig Society* (eds N. Platnick and V. Funk. New York, Columbia University Press. Abridged and reprinted in E. Sober (ed.) pp. 7–36. *Conceptual Issues in Evolutionary Biology*, Cambridge, MA, MIT Press, 1994 (page references to the latter).
- Farris, J.S. (1986). On the boundaries of phylogenetic systematics. *Cladistics* **2**: 14–27.
- Farris, J.S. (1988). *Hennig86*. Published by the author, Port Jefferson Station, New York.
- Farris, J.S. (1989a). The retention index and the rescaled consistency index. *Cladistics* **5**: 417–419.
- Farris, J.S. (1989b). Entropy and fruit flies. *Cladistics* **5**: 103–108.
- Farris, J.S. (1991). Hennig defined paraphyly. *Cladistics* **7**: 297–304.
- Farris, J.S. (1997). Cycles. *Cladistics* **13**: 131–144.
- Farris, J.S. (1999). Likelihood and inconsistency. *Cladistics* **15**: 199–204.
- Farris, J.S. (2001). Support weighting. *Cladistics* **17**: 389–394.

- Farris, J.S. and Kluge, A.G. (1986). Synapomorphy, parsimony, and evidence. *Taxon* **35**: 298–315.
- Farris, J.S., Kluge, A.G. and Eckardt, M.J. (1970). A numerical approach to phylogenetic systematics. *Syst. Zool.* **19**: 172–189.
- Farris, J.S., Källersjö, M., Albert, V.A., Allard, M., Anderberg, A., Bowditch, B., Bult, C., Carpenter, J.M., Crowe, T.M., De Laet, J. *et al.* (1995). Explanation. *Cladistics* **11**: 211–218.
- Farris, J.S., Albert, V.A., Källersjö, M., Lipscomb, D. and Kluge, A.G. (1996). Parsimony jackknifing outperforms neighbor-joining. *Cladistics* **12**: 99–124.
- Farris, J.S., Källersjö, M. and De Laet, J.E. (2001a). Branch lengths do not indicate support — even in maximum likelihood. *Cladistics* **17**: 298–299.
- Farris, J.S., Kluge, A.G. and De Laet, J.E. (2001b). Taxic revisions. *Cladistics* **17**: 79–103.
- Fedorov, A., Merican, A.F. and Gilbert, W. (2002). Large-scale comparison of intron positions among animal, plant, and fungal genes. *Proc. Natl. Acad. Sci. USA* **99**: 16128–16133.
- Felsenstein, J. (1968). *Statistical Inference and the Estimation of Phylogenies*. PhD Thesis, University of Chicago, Chicago, IL.
- Felsenstein, J. (1973). Maximum likelihood and minimum-steps methods for estimating evolutionary trees from data on discrete characters. *Syst. Zool.* **22**: 240–249.
- Felsenstein, J. (1978a). Cases in which parsimony and compatibility methods can be positively misleading. *Syst. Zool.* **27**: 401–410.
- Felsenstein, J. (1978b). The number of evolutionary trees. *Syst. Zool.* **27**: 27–33.
- Felsenstein, J. (1979). Alternative methods of phylogenetic inference and their interrelationship. *Syst. Zool.* **28**: 49–62.
- Felsenstein, J. (1981a). Evolutionary trees from DNA sequences: A maximum likelihood approach. *J. Mol. Evol.* **17**: 368–376.
- Felsenstein, J. (1981b). Evolutionary trees from gene frequencies and quantitative characters: Finding maximum likelihood estimates. *Evolution* **35**: 1229–1242.
- Felsenstein, J. (1981c). A likelihood approach to character weighting and what it tells us about parsimony and compatibility. *Biol. J. Linn. Soc.* **16**: 183–196.
- Felsenstein, J. (1982). Numerical methods for inferring evolutionary trees. *Q. Rev. Biol.* **57**: 379–404.
- Felsenstein, J. (1983). Methods for inferring phylogenies: A statistical view. In *Numerical Taxonomy* (ed. J. Felsenstein), pp. 315–334. Berlin, Springer-Verlag.
- Felsenstein, J. (1988). Phylogenies from molecular sequences: Inference and reliability. *Annu. Rev. Genet.* **2**: 521–565.
- Felsenstein, J. (1996). Inferring phylogenies from protein sequences by parsimony, distance, and likelihood methods. *Methods Enzymol.* **266**: 418–427.
- Felsenstein, J. (2004). *Inferring Phylogenies*. Sunderland, MA, Sinauer Associates.
- Felsenstein, J. and Sober, E. (1987). Parsimony and likelihood: An exchange. *Syst. Zool.* **35**: 617–626.
- Feng, D.F. and Doolittle, R.F. (1987). Progressive sequence alignment as a prerequisite to correct phylogenetic trees. *J. Mol. Evol.* **25**: 351–360.
- Fink, W.L. (1982). The conceptual relationship between ontogeny and phylogeny. *Paleobiology* **8**: 254–264.
- Fitch, W.M. (1970). Distinguishing homologous from analogous proteins. *Syst. Zool.* **19**: 99–113.
- Fitch, W.M. (1971). Toward defining the course of evolution: Minimal change for a specific tree topology. *Syst. Zool.* **20**: 406–416.
- Fitz-Gibbon, S.T. and House, C.H. (1999). Whole genome-based phylogenetic analysis of free-living microorganisms. *Nucleic Acids Res.* **27**: 4218–4222.
- Force, A., Lynch, M., Pickett, F.B., Amores, A., Yan, Y.L. and Postlethwait, J. (1999). Preservation of duplicate genes by complementary, degenerative mutations. *Genetics* **151**: 1531–1545.
- Forster, M.R. (2000). Key concepts in model selection: Performance and generalizability. *J. Math. Psych.* **44**: 205–231.
- Forster, M.R. and Sober, E. (1994). How to tell when simpler, more unified, or less *ad hoc* theories will provide more accurate predictions. *Br. J. Phil. Sci.* **45**: 1–35.
- Foulds, L.R. (1984). Maximum savings in the Steiner problem in phylogeny. *J. Theoret. Biol.* **107**: 471–474.
- Foulds, L.R. and Graham, R.L. (1982). The Steiner problem in phylogeny is NP-complete. *Adv. Appl. Math.* **3**: 43–49.
- Friedman, M. (1983). *Foundations of Space-Time Theories: Relativistic Physics and Philosophy of Science*. Princeton, NJ, Princeton University Press.
- Fredman, M.L. (1984). Algorithms for computing evolutionary similarity measures with length independent gap penalties. *Bull. Math. Biol.* **46**: 545–563.
- Freudenstein, J.V., Pickett, K.M., Simmons, M.P. and Wenzel, J.W. (2003). From basepairs to birdsongs: phylogenetic data in the age of genomics. *Cladistics* **19**: 333–347.
- Frost, D.R. (2000). Species, descriptive efficiency, and progress in systematics. In *The Biology of Plethodontid Salamanders* (eds R.C. Bruce, R.J. Jaeger, and L.D. Houck), pp. 7–29. New York, Kluwer Academic/Plenum Publishing.

- Frost, D.R. and Kluge, A.G. (1994). A consideration of epistemology in systematic biology, with special reference to species. *Cladistics* **10**: 259–294.
- Frost, D.R., Rodrigues, M.T., Grant, T. and Titus, T.A. (2001). Phylogenetics of the lizard genus *Tropidurus* (Squamata: Tropiduridae: Tropidurinae): Direct optimization, descriptive efficiency, and sensitivity analysis of congruence between molecular data and morphology. *Mol. Phylogenet. Evol.* **21**: 352–371.
- Fukami-Kobayashi, K., Schreiber, D.R. and Benner, S.A. (2002). Detecting compensatory covariation signals in protein evolution using reconstructed ancestral sequences. *J. Mol. Biol.* **319**: 729–743.
- Funk, V.A. and Brooks, D.R. (1990). *Phylogenetic Systematics as the Basis of Comparative Biology*. Washington, DC, Smithsonian Institution Press.
- Gaasterland, T. and Ragan, M.A. (1998). Microbial genescapes: Phyletic and functional patterns of ORF distribution among prokaryotes. *Microb. Comp. Genomics* **3**: 199–217.
- Gallut, C. and Barriol, V. (2002). Cladistic coding of genomic maps. *Cladistics* **18**: 526–536.
- Galperin, M.Y. and Koonin, E.V. (2000). Who's your neighbor? New computational approaches for functional genomics. *Nat. Biotechnol.* **18**: 609–613.
- Gee, H. (2000). *Deep Time: Cladistics, the Revolution in Evolution*. London, Fourth Estate.
- Ghiselin, M.T. (1966). On semantic pitfalls of biological adaptation. *Philos. Sci.* **33**: 147–153.
- Ghiselin, M.T. (2004). Mayr and Bock versus Darwin on genealogical classification. *J. Zool. Syst. Evol. Res.* **42**: 165–169.
- Giribet, G. (2002). Relationships among metazoan phyla as inferred from 18S rRNA sequence data: A methodological approach. In *Molecular Systematics and Evolution: Theory and Practice*, (eds R. DeSalle, G. Giribet, and W. Wheeler), pp. 85–101. Basel, Birkhäuser Verlag.
- Giribet, G. and Wheeler, G. (1999). On gaps. *Mol. Phylogenet. Evol.* **13**: 132–143.
- Giribet, G., Distel, D.L., Polz, M., Sterrer, W. and Wheeler, W.C. (2000). Triploblastic relationships with emphasis on the acoelomates and the position of Gnathostomulida, Cyclophora, Plathelminthes, and Chaetognatha: A combined approach of 18S rDNA sequences and morphology. *Syst. Biol.* **49**: 539–562.
- Giribet, G., Wheeler, W.C. and Muona, J. (2002). DNA multiple sequence alignments. In *Molecular Systematics and Evolution: Theory and Practice* (eds R. Desalle, G. Giribet and W. Wheeler), pp. 107–114. Basel, Birkhäuser Verlag.
- Gladstein, D.S. (1997). Efficient incremental character optimization. *Cladistics* **13**: 21–26.
- Gogarten, J.P., Doolittle, W.F. and Lawrence, J.G. (2002). Prokaryotic evolution in light of gene transfer. *Mol. Biol. Evol.* **19**: 2226–2238.
- Goldman, N. (1990). Maximum likelihood inference of phylogenetic trees, with special reference to a Poisson process model of DNA substitution and to parsimony analyses. *Syst. Zool.* **39**: 345–361.
- Goloboff, P.A. (1995). A revision of the south American spiders of the family Nemesiidae (Araneae, Mygalomorphae). Part I: species from Peru, Chile, Argentina, and Uruguay. *Bull. Am. Mus. Nat. Hist.* **224**: 1–189.
- Goloboff, P.A. (1993a). Estimating character weights during tree search. *Cladistics* **9**: 83–91.
- Goloboff, P.A. (1993b). *Nona*: a tree-searching program. Available at <http://www.zmuc.dk/public/phylogeny/Nona-PeeWee/>.
- Goloboff, P.A. (1993c). *Pee-Wee*: Parsimony and Implied weights. Available at <http://www.zmuc.dk/public/phylogeny/Nona-PeeWee/>.
- Goloboff, P.A. (1994). Character optimization and calculation of tree lengths. *Cladistics* **9**: 433–436.
- Goloboff, P.A. (1995). *SPA: Sankoff Parsimony Analysis*, ver. 1.1. Available at <http://www.zmuc.dk/public/phylogeny/Nona-PeeWee/>.
- Goloboff, P.A. (1996a). *PHAST: Phylogenetic Analysis for Sankovian Transformations*, ver. 1.1. Available at <http://www.zmuc.dk/public/phylogeny/Nona-PeeWee/>.
- Goloboff, P.A. (1996b). Methods for faster parsimony analysis. *Cladistics* **12**: 199–220.
- Goloboff, P.A. (1998b). Tree searches under Sankoff parsimony. *Cladistics* **14**: 229–238.
- Goloboff, P.A. (1999). Analyzing large data sets in reasonable times: Solutions for composite optima. *Cladistics* **15**: 415–428.
- Goloboff, P.A. (2003). Parsimony, likelihood, and simplicity. *Cladistics* **19**: 91–103.
- Goloboff, P.A. and Farris, J.S. (2001). Methods for quick consensus estimation. *Cladistics* **17**: S26–S34.
- Goloboff, P.A., Wheeler, W. and Pol, D. (2003a). Parallel TNT. *Cladistics* **19**: 152 (in Muona, J. (2003). Abstracts of the 21st annual meeting of the Willi Hennig society. *Cladistics* **19**: 148–163.)
- Goloboff, P., Farris, J., Källersjö, M., Oxelmann, B., Ramírez, M. and Szumik, C. (2003b). Improvements to resampling measures of group support. *Cladistics* **19**: 324–332.
- Goloboff, P., Farris, J. and Nixon, K. (2004). *T.N.T.: Tree Analysis Using New Technology*. Available at www.zmuc.dk/public/phylogeny/tnt.
- Goodman, M., Czelusniak, J., Moore, G.W., Romero-Herrera, A.E. and Matsuda, G. (1979). Fitting the gene

- lineage into its species lineage, a parsimony strategy illustrated by cladograms constructed from globin sequences. *Syst. Zool.* **28**: 132–163.
- Gouldge, T.A. (1961). *The Ascent of Life. A Philosophical Study of the Theory of Evolution*. Toronto, University of Toronto Press. [1967 reprint]
- Grant, T. (2002). Testing methods: The evaluation of discovery operations in evolutionary biology. *Cladistics* **18**: 94–111.
- Grant, T. and Kluge, A.G. (2003). Data exploration in phylogenetic inference: Scientific, heuristic, or neither. *Cladistics* **19**: 379–418.
- Grant, T. and Kluge, A.G. (2004). Transformation series as an ideographic character concept. *Cladistics* **20**: 23–31.
- Greene, B. (2004). *The Fabric of the Cosmos. Space, Time, and the Texture of Reality*. New York, A.A. Knopf.
- Greuter, W., McNeill, J., Barrie, F.R., Burdet, H.M., Demoulin, V., Filgueiras, T.S., Nicolson, D.H., Silva, P.C., Skog, J.E., Trehane, P., et al. (2000). *International Code of Botanical Nomenclature (St. Louis Code)*. *Regnum Vegetabile* 138. Königstein, Koeltz Scientific Books.
- Gu, J. and Gu, X. (2003). Induced gene expression in human brain after the split from chimpanzee. *Trends Genet.* **19**: 63–65.
- Gusfield, D. (1997). *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*. Cambridge, Cambridge University Press.
- Hacking, I. (1965). *The Logic of Statistical Inference*. Cambridge, Cambridge University Press.
- Hartigan, J.A. (1973). Minimum mutation fits to a given tree. *Biometrics* **29**: 53–65.
- Harvey, P.H. and Pagel, M.D. (1991). *The Comparative Method in Evolutionary Biology*. New York, Oxford University Press.
- Hasegawa, M. and Kishino, H. (1989). Confidence limits on the maximum-likelihood estimate of the hominoid tree from mitochondrial-DNA sequences. *Evolution* **43**: 672–677.
- Hastings, W.K. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* **57**: 97–109.
- Hedges, S.B. (2002). The origin and evolution of model organisms. *Nat. Rev. Genet.* **3**: 838–849.
- Hein, J. (1989a). A new method that simultaneously aligns and reconstructs ancestral sequences for any number of homologous sequences when a phylogeny is given. *Mol. Biol. Evol.* **6**: 649–668.
- Hein, J. (1989b). A tree reconstruction method that is economical in the number of pairwise comparisons used. *Mol. Biol. Evol.* **6**: 669–684.
- Hein, J.J. (2001). An algorithm for statistical alignment of sequences related by a binary tree. In *Pacific Symposium on Biocomputing 2001*. (eds R.B. Altman, A.K. Dunker, L. Hunter, K. Lauderdale and T.E. Klein), vol. 6, pp. 179–190. Singapore, World Scientific.
- Hein, J., Jensen, J.L. and Pedersen, C.N.S. (2003). Recursions for statistical multiple alignment. *Proc. Natl. Acad. Sci. USA* **100**: 14960–14965.
- Hendy, M.D. and Penny, D. (1982). Branch and bound algorithms to determine minimal evolutionary trees. *Math. Biosci.* **59**: 277–290.
- Hendy, M.D., Foulds, L.R. and Penny, D. (1980). Proving phylogenetic trees minimal with I-clustering and set partitioning. *Math. Biosci.* **51**: 71–88.
- Hennig, W. (1950). *Grundzüge einer Theorie der phylogenetischen Systematik*. Berlin, Deutscher Zentralverlag.
- Hennig, W. (1966). *Phylogenetic Systematics*. Urbana, IL, University of Illinois Press.
- Higgins, D.G. and Sharp, P.M. (1988). CLUSTAL: A package for performing multiple sequence alignment on a microcomputer. *Gene* **73**: 237–244.
- Huber, K.T., Moulton, V. and Steel, M. (2002). *Four characters suffice to convexly define a phylogenetic tree*. Research Report UCDMA2002/12, Christchurch, New Zealand, Department of Mathematics and Statistics, University of Canterbury.
- Huelsenbeck, J.P. and Lander, K.M. (2003). Frequent inconsistency of parsimony under a simple model of cladogenesis. *Syst. Biol.* **52**: 641–648.
- Huelsenbeck, J.P. and Ronquist, F. (2001). MrBayes: Bayesian inference of phylogeny. *Bioinformatics* **17**: 754–755.
- Huelsenbeck, J.P., Bull, J.J. and Cunningham, C.W. (1996). Combining data in phylogenetic analysis. *Trends Ecol. Evol.* **4**: 152–158.
- Huelsenbeck, J., Ronquist, F., Nielsen, R. and Bollback, J. (2001). Bayesian inference of phylogeny and its impact on evolutionary biology. *Science* **294**: 2310–2314.
- Huelsenbeck, J.P., Larget, B., Miller, R.E. and Ronquist, F. (2002). Potential applications and pitfalls of Bayesian inference of phylogeny. *Syst. Biol.* **51**: 673–688.
- Huelsenbeck, J., Larget, B. and Alfaro, M. (2004). Bayesian phylogenetic model selection using reversible jump Markov chain Monte Carlo. *Mol. Biol. Evol.* **21**: 1123–1133.
- Hull, D.L. (1967). Certainty and circularity in evolutionary taxonomy. *Evolution* **21**: 174–189.
- Hull, D.L. (1974). *Philosophy of Biological Science*. Englewood Cliffs, NJ, Prentice-Hall.
- Hull, D.L. (1975). Central subjects and historical narratives. *History and Theory: Studies Philos. Hist.* **14**: 253–274.
- Hull, D.L. (1977). The ontological status of species as evolutionary units. In *Foundational Problems in Special Sciences* (eds R. Butts and J. Hintikka), pp. 91–102. Dordrecht, D. Reidel Pub. Co.

- Hull, D.L. (1981). Historical narratives and integrating explanations. In *Pragmatism and Purpose: Essays Presented to Thomas A. Goudge* (eds L.W. Sumner, J.G. Slater and F. Wilson), pp. 172–188, 308–310. Toronto, University of Toronto Press.
- Hull, D.L. (1982). Exemplars and scientific change. *Proc. Biennial Mtg. Phil. Sci. Assoc.* **2**: 479–503.
- Hull, D.L. (1989). *The Metaphysics of Evolution*. Albany, NY, SUNY Press.
- Huson, D.H. and Steel, M. (2004). Phylogenetic trees based on gene content. *Bioinformatics* **20**: 2044–2049.
- Huynen, M.A. and Bork, P. (1998). Measuring genome evolution. *Proc. Natl. Acad. Sci. USA* **95**: 5849–5856.
- ICZN (1999). *International Code of Zoological Nomenclature, 4th Edn.* London, International Trust for Zoological Nomenclature.
- Jacob, F., Perrin, D., Sanchez, C. and Monod, J. (1960). L'Operon: groupe de genes a expression coordonnee par un operateur. *C. R. Seance Acad. Sci.* **250**: 1727–1729.
- Jenner, R.A. (2004). The scientific status of metazoan cladistics: why current reserach practice must change. *Zool. Scripta* **33**: 293–310.
- Jermann, T.M., Opitz, J.G., Stackhouse, J. and Benner, S.A. (1995). Reconstructing the evolutionary history of the artiodactyl ribonuclease superfamily. *Nature* **374**: 57–59.
- Jiang, T.L. and Lawler, E.L. (1994). Aligning sequences via an evolutionary tree: Computational complexity and approximation. In *Proceedings of the 26th ACM Symposium on the Theory of Computing*, pp. 760–769. New York, ACM.
- Källersjö, M., Farris, J.S., Chase, M.W., Bremer, B., Fay, M.F., Humphries, C.J., Petersen, G., Seberg, O. and Bremer, K. (1998). Simultaneous parsimony jackknife analysis of 2538 *rbcl* DNA sequences reveals support for major clades of green plants, land plants, seed plants and flowering plants. *Plant Syst. Evol.* **213**: 259–287.
- Källersjö, M., Albert, V.A. and Farris, J.S. (1999). Homoplasy increases phylogenetic structure. *Cladistics* **15**: 91–93.
- Kapitonov, V.V. and Jurka, J. (1999). The long terminal repeat of an endogenous retrovirus induces alternative splicing and encodes an additional carboxy-terminal sequence in the human leptin receptor. *J. Mol. Evol.* **48**: 248–251.
- Katinka, M.D., Duprat, S., Cornillot, E., Metenier, G., Thomarat, F., Prensier, G., Barbe, V., Peyretailade, E., Brottier, P., Wincker, P. *et al.* (2001). Genome sequence and gene compaction of the eukaryote parasite *Encephalitozoon cuniculi*. *Nature* **414**: 450–453.
- Kidd, K.K. and Sgaramella-Zonta, L.A. (1971). Phylogenetic analysis: Concepts and methods. *Am. J. Hum. Genet.* **23**: 235–252.
- Kim, J. (1996). General inconsistency conditions for maximum parsimony: Effects of branch lengths and increasing numbers of taxa. *Syst. Biol.* **45**: 363–374.
- Kimura, M. and Crow, J. (1964). The number of alleles that can be maintained in a finite population. *Genetics* **49**: 725–738.
- Kjer, K.M. (1995). Use of rRNA secondary structure in phylogenetic studies to identify homologous positions: An example of alignment and data presentation from the frogs. *Mol. Phylogenet. Evol.* **4**: 314–330.
- Kluge, A.G. (1988). Parsimony in vicariance biogeography: A quantitative method and a Greater Antillean example. *Syst. Zool.* **37**: 315–328.
- Kluge, A.G. (1989). A concern for evidence and a phylogenetic hypothesis of relationships among *Epicrates* (Boidae, Serpentes). *Syst. Zool.* **38**: 7–25.
- Kluge, A.G. (1997a). Testability and the refutation and corroboration of cladistic hypotheses. *Cladistics* **13**: 81–96.
- Kluge, A.G. (1997b). Sophisticated falsification and research cycles: Consequences for differential character weighting in phylogenetic systematics. *Zool. Scripta* **26**: 349–360.
- Kluge, A.G. (1999). The science of phylogenetic systematics: Explanation, prediction, and test. *Cladistics* **15**: 429–436.
- Kluge, A.G. (2001a). Parsimony with and without scientific justification. *Cladistics* **17**: 199–210.
- Kluge, A.G. (2001b). Philosophical conjectures and their refutation. *Syst. Biol.* **50**: 322–330.
- Kluge, A.G. (2002). Distinguishing “or” from “and” and the case for historical identification. *Cladistics* **18**: 585–593.
- Kluge, A.G. (2003a). On the deduction of species relationships: A précis. *Cladistics* **19**: 233–239.
- Kluge, A.G. (2003b). The repugnant and the mature in phylogenetic inference: A temporal similarity and historical identity. *Cladistics* **19**: 356–368.
- Kluge, A.G. (2004). On total evidence: For the record. *Cladistics* **20**: 205–207.
- Kluge, A.G. (2005). Taxonomy in theory and practice, with arguments for a new phylogenetic system of taxonomy. In *Ecology and Evolution in the Tropics: a Herpetological Perspective* (eds M.A. Donnelly, B.I. Crother, C. Guyer, M.H. Wake and M.E. White), pp. 7–47. Chicago, University of Chicago Press.
- Kluge, A.G. and Farris, J.S. (1969). Quantitative phyletics and the evolution of Anurans. *Syst. Zool.* **18**: 1–32.
- Kluge, A.G. and Farris, J.S. (1999). Taxic homology = overall similarity. *Cladistics* **15**: 205–212.
- Koonin, E.V., Aravind, L. and Kondrashov, A.S. (2000). The impact of comparative genomics on our understanding of evolution. *Cell* **101**: 573–576.

- Koonin, E.V., Makarova, K.S. and Aravind, L. (2001). Horizontal gene transfer in prokaryotes: Quantification and classification. *Annu. Rev. Microbiol.* **55**: 709–742.
- Koonin, E.V. and Galperin, M.Y. (2002). *Sequence-Evolution-Function. Computational Approaches in Comparative Genomics*. Kluwer Academic Publishers, New York.
- Koonin, E.V., Fedorova, N.D., Jackson, J.D., Jacobs, A.R., Krylov, D.M., Makarova, K.S., Mazumder, R., Mekhedov, S.L., Nikolskaya, A.N., Rao, B.S., *et al.* (2004). A comprehensive evolutionary classification of proteins encoded in complete eukaryotic genomes. *Genome Biol.* **5**: R7.
- Korbel, J.O., Snel, B., Huynen, M.A. and Bork, P. (2002). SHOT: A web server for the construction of genome phylogenies. *Trends Genet.* **18**: 158–162.
- Koshi, J.M. and Goldstein, R.A. (1996). Probabilistic reconstruction of ancestral protein sequences. *J. Mol. Evol.* **42**: 313–320.
- Kruskal, J. (1983). An overview of sequence comparison. In *Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison* (eds D. Sankoff and J. Kruskal), pp. 1–44. Stanford, CA, CSLI Publications (1999 reprint).
- Kumar, S., Tamura, K. and Nei, M. (1993). MEGA: *Molecular Evolutionary Genetics Analysis, vers. 1.01*. University Park, PA, Pennsylvania State University.
- Kunin, V. and Ouzounis, C.A. (2003). The balance of driving forces during genome evolution in prokaryotes. *Genome Res.* **13**: 1589–1594.
- Lande, R. (1976). Natural selection and random genetic drift in phenotypic evolution. *Evolution* **30**: 314–334.
- Larget, B. and Simon, D. (1999). Markov chain Monte Carlo algorithms for the Bayesian analysis of phylogenetic trees. *Mol. Biol. Evol.* **16**: 750–759.
- Larson, A. and Losos, J.B. (1996). Phylogenetic systematics of adaptation. In *Adaptation* (eds M.R. Rose and G.V. Lauder), pp. 187–220. San Diego, CA, Academic Press.
- Laudan, L. (1990). *Science and Relativism: Some Key Controversies in the Philosophy of Science*. Chicago, University of Chicago Press.
- Laudan, R. (1990). What's so special about the past? In *Evolutionary Innovations* (ed. M. Nitechi), pp. 55–67. Chicago, University of Chicago Press.
- Lauder, G.V., Leroi, A.M. and Rose, M.R. (1993). Adaptations and History. *Trends Ecol. Evol.* **8**: 294–297.
- Lawrence, J. (1999). Selfish operons: The evolutionary impact of gene clustering in prokaryotes and eukaryotes. *Curr. Opin. Genet. Dev.* **9**: 642–648.
- Lawrence, J.G. and Roth, J.R. (1996). Selfish operons: Horizontal transfer may drive the evolution of gene clusters. *Genetics* **143**: 1843–1860.
- Le Cam, L. (1960). An approximation theorem for the Poisson binomial distribution. *Pacific J. Math.* **10**: 1181–1197.
- Le Quesne, W. (1969). A method of selection of characters in numerical taxonomy. *Syst. Zool.* **18**: 201–205.
- Lee, D.C. and Bryant, H.N. (1999). A reconsideration of the coding of inapplicable characters: Assumptions and problems. *Cladistics* **15**: 373–378.
- Levesque, M., Shasha, D., Kim, W., Surette, M.G. and Benfey, P.N. (2003). Trait-to-gene: a computational method for predicting the function of uncharacterized genes. *Curr. Biol.* **13**: 129–133.
- Lewis, P.O. (2001). A likelihood approach to estimating phylogeny from discrete morphological character data. *Syst. Biol.* **50**: 913–925.
- Li, S. (1996). *Phylogenetic Tree Construction using Markov Chain Monte Carlo*. PhD Dissertation, Ohio State University, Columbus, OH.
- Li, S., Pearl, D.K. and Doss, H. (2000). Phylogenetic tree construction using Markov chain Monte Carlo. *J. Am. Stat. Assoc.* **2000**: 493–508.
- Liberles, D.A. (2001). Evaluation of methods for determination of a reconstructed history of gene sequence evolution. *Mol. Biol. Evol.* **18**: 2040–2047.
- Liberles, D.A., Schreiber, D.R., Govindarajan, S., Chamberlin, S.G. and Benner, S.A. (2001). The Adaptive Evolution Database (TAED). *Genome Biol.* **2**(8): research0028.1– research0028.6.
- Liberles, D.A., Thoren, A., von Heijne, G. and Elofsson, A. (2002). The use of phylogenetic profiles for gene predictions. *Curr. Genomics* **3**: 131–137.
- Lidén, M. (1990). Replicators, hierarchy, and the species problem. *Cladistics* **6**: 183–186.
- Lipscomb, D.L. (1992). Parsimony, homology, and the analysis of multistate characters. *Cladistics* **8**: 45–65.
- Logsdon, Jr., J.M., Stoltzfus, A. and Doolittle, W.F. (1998). Molecular evolution: Recent cases of spliceosomal intron gain? *Curr. Biol.* **8**: R560–R63.
- Logsdon, Jr., J.M., Tyshenko, M.G., Dixon, C., J, D.J., Walker, V.K. and Palmer, J.D. (1995). Seven newly discovered intron positions in the triose-phosphate isomerase gene: Evidence for the introns-late theory. *Proc. Natl. Acad. Sci. USA* **92**: 8507–8511.
- Long, M. and Rosenberg, C. (2000). Testing the “proto-splice sites” model of intron origin: Evidence from analysis of intron phase correlations. *Mol. Biol. Evol.* **17**: 1789–1796.
- Lutzoni, F., Wagner, P., Reeb, V. and Zoller, S. (2000). Integrating ambiguously aligned regions of DNA sequences in phylogenetic analyses without violating positional homology. *Syst. Biol.* **49**: 628–651.
- Lynch, M. and Richardson, A.O. (2002). The evolution of spliceosomal introns. *Curr. Opin. Genet. Dev.* **12**: 701–710.

- Maddison, D.R. (1991). The discovery and importance of multiple islands of most-parsimonious trees. *Syst. Zool.* **40**: 315–328.
- Maddison, D.R. and Maddison, W.P. (2001). *MacClade 4: Analysis of Phylogeny and Character Evolution* (incl. vers. 4.03). Sunderland, MA, Sinauer Associates.
- Maddison, D.R., Swofford, D.L. and Maddison, W.P. (1997). NEXUS: An extensible file format for systematic information. *Syst. Biol.* **46**: 590–621.
- Maddison, W.P. (1991). Squared-change parsimony reconstructions of ancestral states for continuous-valued characters on a phylogenetic tree. *Syst. Zool.* **40**: 304–314.
- Maddison, W.P. (1993). Missing data versus missing characters in phylogenetic analysis. *Syst. Biol.* **42**: 576–581.
- Maddison, W.P. and Maddison, D.R. (1992). *MacClade 3: Analysis of Phylogeny and Character Evolution* (incl. vers. 3.04). Sunderland, MA, Sinauer Associates.
- Mallatt, J. and Winchell, C.J. (2002). Testing the new animal phylogeny: First use of combined large-subunit and small-subunit rRNA gene sequences to classify the protostomes. *Mol. Biol. Evol.* **19**: 289–301.
- Marchionni, M. and Gilbert, W. (1986). The triosephosphate isomerase gene from maize: Introns antedate the plant-animal divergence. *Cell* **46**: 133–141.
- Marcotte, E.M., Pellegrini, M., Thompson, M.J., Yeates, T.O. and Eisenberg, D.A. (1999). A combined algorithm for genome-wide prediction of protein function. *Nature* **402**: 83–86.
- Martin C. and Paz-Ares J. (1997). MYB transcription factors in plants. *Trends Plant Sci.* **13**: 67–73.
- Mau, B. (1996). *Bayesian Phylogenetic Inference via Markov Chain Monte Carlo Methods*. PhD Dissertation, University of Wisconsin, Madison, WI.
- Mau, B. and Newton, M. (1997). Phylogenetic inference for binary data on dendrograms using Markov chain Monte Carlo. *J. Comput. Graph. Stat.* **6**: 122–131.
- Mau, B., Newton, M. andarget, B. (1999). Bayesian phylogenetic inference via Markov chain Monte Carlo methods. *Biometrics* **55**: 1–12.
- Mayr, E. and Bock, W.J. (2002). Classification and other ordering systems. *J. Zool. Syst. Evol. Res.* **40**: 169–194.
- McAllister, J.W. (1996). *Beauty and Revolution in Science*. Ithaca, NY, Cornell University Press.
- McAllister, J.W. (2000). Unification of theories. In *A Companion to the Philosophy of Science* (ed. W.H. Newton-Smith), pp. 537–539. Oxford, Blackwell Publishing.
- McDade, L.A. (1992). Hybrids and phylogenetic systematics II. The impact of hybrids on cladistic analysis. *Evolution* **46**: 1329–1346.
- Messier, W. and Stewart, C.B. (1997). Episodic adaptive evolution of primate lysozymes. *Nature* **385**: 151–154.
- Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H. and Teller, E. (1953). Equations of state calculations by fast computing machines. *J. Chem. Phys.* **21**: 1087–1091.
- Mickevich, M.F. and Farris, J.S. (1981). *PHYSYS: Phylogenetic Analysis System*. Published by the authors.
- Miklós, I., Lunter, A. and Holmes, I. (2004). A “long indel” model for evolutionary sequence alignment. *Mol. Biol. Evol.* **21**: 529–540.
- Mindell, D.P. and Thacker, C.E. (1996). Rates of molecular evolution: Phylogenetic issues and applications. *Annu. Rev. Ecol. Syst.* **27**: 279–303.
- Mirkin, B.G., Fenner, T.I., Galperin, M.Y. and Koonin, E.V. (2003). Algorithms for computing parsimonious evolutionary scenarios for genome evolution, the last universal common ancestor and dominance of horizontal gene transfer in the evolution of prokaryotes. *BMC Evol. Biol.* **3**: 2.
- Mishler, B.D. (1994). Cladistic analysis of molecular and morphological data. *Am. J. Phys. Anthropol.* **94**: 143–156.
- Mishler, B.D. (1999). Getting rid of species? In *Species: New Interdisciplinary Essays* (ed. R. Wilson), pp. 307–315. Cambridge, MA, MIT Press.
- Mishler, B.D. (2000). Deep phylogenetic relationships among “plants” and their implications for classification. *Taxon* **49**: 661–683.
- Mishler, B.D. and Brandon, R.N. (1987). Individuality, pluralism, and the phylogenetic species concept. *Biol. Phil.* **2**: 397–414.
- Mishler, B.D. and De Luna, E. (1991). The use of ontogenetic data in phylogenetic analyses of mosses. *Adv. Bryol.* **4**: 121–167.
- Mishler, B.D. and Theriot, E. (2000a). The phylogenetic species concept (*sensu* Mishler and Theriot): Monophyly, apomorphy, and phylogenetic species concepts. In *Species Concepts and Phylogenetic Theory: A Debate*. (eds Q.D. Wheeler and R. Meier), pp. 44–54. New York, Columbia University Press.
- Mishler, B.D. and Theriot, E.G. (2000b). A critique from the Mishler and Theriot phylogenetic species concept perspective: Monophyly, apomorphy, and phylogenetic species concepts. In *Species Concepts and Phylogenetic Theory: A Debate* (eds Q.D. Wheeler and R. Meier), pp. 133–145. New York, Columbia University Press.
- Mishler, B.D. and Theriot, E.G. (2000c). A defense of phylogenetic species concept (*sensu* Mishler and Theriot): Monophyly, apomorphy, and phylogenetic species concepts. In *Species Concepts and Phylogenetic Theory: A Debate* (eds Q.D. Wheeler and R. Meier), pp. 179–184. New York, Columbia University Press.
- Modrek, B. and Lee, C.J. (2003). Alternative splicing in the human, mouse, and rat genomes is associated with

- an increased frequency of exon creation and/or loss. *Nat. Genet.* **34**: 177–180.
- Moilanen, A. (1999). Searching for most parsimonious trees with simulated evolutionary optimization. *Cladistics* **15**: 39–50.
- Moilanen, A. (2001). Simulated evolutionary optimization and local search: Introduction and application to tree search. *Cladistics* **17**: S12–S25.
- Montague, M.G. and Hutchison, 3rd., C.A. (2000). Gene content phylogeny of herpesviruses. *Proc. Natl. Acad. Sci. USA* **97**: 5334–5339.
- Moran, N.A. (2002). Microbial minimalism: Genome reduction in bacterial pathogens. *Cell* **108**: 583–586.
- Moret, B.M.E., Wang, L.S., Warnow, T. and Wyman, S. (2001). New approaches for reconstructing phylogenies based on gene order. Proc. 9th Int'l Conf. on Intelligent Systems for Molecular Biology ISMB-2001, *Bioinformatics* **17**: S165–S173.
- Moret, B.M.E., Tang, J., Wang, L.S. and Warnow, T. (2002). Steps toward accurate reconstruction of phylogenies from gene-order data. *J. Comput. Syst. Sci.* **65**(3): 508–525.
- Morgenstern, B. (2004). DIALIGN: Multiple DNA and protein sequence alignment at BiBiServ. *Nucleic Acids Res.* **32**: W33–W36.
- Morgenstern, B., Dress, A. and Werner, T. (1996). Multiple DNA and protein sequence alignment based on segment-to-segment comparison. *Proc. Natl. Acad. Sci. USA* **93**: 12098–12103.
- Moritz, C. (2002). Strategies to protect biological diversity and the processes that sustain it. *Syst. Biol.* **51**: 238–254.
- Mossel, E. and Steel, M. (2004a). A phase transition for a random cluster model on phylogenetic trees. *Math. Biosci.* **187**: 189–203.
- Mossel, E. and Steel, M. (2005). How much can evolved characters tell us about the tree that generated them? In *Mathematics of Evolution and Phylogeny* (ed. O. Gascuel). Oxford, Oxford University Press.
- Murata, M., Richardson, J. and Sussman, J. (1985). Simultaneous comparison of three protein sequences. *Proc. Natl. Acad. Sci. USA* **82**: 3073–3077.
- Mushegian, A.R. and Koonin, E.V. (1996). Gene order is not conserved in bacterial evolution. *Trends Genet.* **12**: 289–290.
- Mushegian, A.R., Garey, J.R., Martin, J. and Liu, L.X. (1998). Large-scale taxonomic profiling of eukaryotic model organisms: A comparison of orthologous proteins encoded by the human, fly, nematode, and yeast genomes. *Genome Res.* **8**: 590–598.
- Myllykallio, H., Lipowski, G., Leduc, D., Filee, J., Forterre, P. and Liebl, U. (2002). An alternative flavin-dependent mechanism for thymidylate synthesis. *Science* **297**: 105–107.
- Nadeau, J.J. and Taylor, B.A. (1984). Lengths of chromosome segments conserved since divergence of man and mouse. *Proc. Natl. Acad. Sci. USA* **81**: 814–818.
- Naylor, G.J.P. and Adams D.C. (2001). Are the fossil data really at odds with the molecular data? Morphological evidence for *Cetartiodactyla* phylogeny reexamined. *Syst. Biol.* **50**: 444–453.
- Naylor, G.J.P. and Adams, D.C. (2003). Total evidence versus relevant evidence: A response to O'Leary *et al.* (2003). *Syst. Biol.* **52**: 864–865.
- Needleman, S.B. and Wunsch, C.D. (1970). A general method applicable to the search for similarities in the amino acid sequence of two proteins. *J. Mol. Biol.* **48**: 443–453.
- Neff, N.A. (1986). A rational basis for *a priori* character weighting. *Syst. Zool.* **35**: 102–109.
- Nei, M. and Kumar, S. (2001). *Molecular Evolution and Phylogenetics*. Oxford, Oxford University press.
- Nelson, G. and Platnick, N.I. (1981). *Systematics and Biogeography: Cladistics and Vicariance*. New York, Columbia University Press.
- Newton, M., Mau, B. and Larget, B. (1999). Markov chain Monte Carlo for the Bayesian analysis of evolutionary trees from aligned molecular sequences. In *Statistics in Molecular Biology and Genetics*, vol. 33 (ed. F. Seillier-Moisewitsch), pp. 143–162. Bethesda, MD, Institute of Mathematical Statistics.
- Neyman, J. (1971). Molecular studies of evolution: A source of novel statistical problems. In *Statistical Decision Theory and Related Topics* (eds S. Gupta and J. Yackel), pp. 1–27. New York, Academic Press.
- Nixon, K.C. (1999). The parsimony ratchet, a new method for rapid parsimony analysis. *Cladistics* **15**: 407–414.
- Nixon, K.C. (2002). *WinClada*, vers. 1.00.08. Published by the author, Ithaca, New York (distributed through www.cladistics.org).
- Nixon, K.C. and Carpenter, J.M. (1993). On outgroups. *Cladistics* **9**: 413–426.
- Nixon, K.C. and Little, D.P. (2004). The use of optimality criteria in DNA sequence data and its application in a new computer program (Abstract, 22nd annual meeting of the Willi Hennig Society, New York 2003). *Cladistics* **20**: 90–91.
- Nixon, K.C. and Wheeler, Q.D. (1990). An amplification of the phylogenetic species concept. *Cladistics* **6**: 211–223.
- Nolan, D. (1997). Quantitative parsimony. *Br. J. Phil. Sci.* **48**: 329–343.

- Notredame, C. (2002). Recent progress in multiple sequence alignment: A survey. *Pharmacogenomics* **3**: 1–14.
- Notredame, C., Holm, L. and Higgins, D.G. (1998). COFFEE: An objective function for multiple sequence alignments. *Bioinformatics* **14**: 407–422.
- Notredame, C., Higgins, D.G. and Heringa, J. (2000). T-Coffee: A novel method for fast and accurate multiple sequence alignment. *J. Mol. Biol.* **302**: 205–217.
- O'Hara, R.J. (1988). Homage to Clio, or, toward an historical philosophy for evolutionary biology. *Syst. Zool.* **37**: 142–155.
- Ochman, H., Lawrence, J.G. and Groisman, E.A. (2000). Lateral gene transfer and the nature of bacterial innovation. *Nature* **405**: 299–304.
- Ochoterena, H. (2004). Independence of alignment and phylogenetic reconstruction and their optimality criteria (Abstract, 22nd annual meeting of the Willi Hennig Society, New York 2003). *Cladistics* **20**: 91.
- Ohno, S. (1970). *Evolution by Gene Duplication*. New York, Springer-Verlag.
- Oleksiak, M.F., Churchill, G.A. and Crawford, D.L. (2002). Variation in gene expression within and among natural populations. *Nat. Genet.* **32**: 261–266.
- Padian, K. (2004). For Darwin, 'genealogy alone' did give classification. *J. Zool. Syst. Evol. Res.* **42**: 162–164.
- Parzen, E. (1962). *Stochastic Processes*. San Francisco, Holden-Day.
- Patterson, C. (1982). Morphological characters and homology. In *Problems of Phylogenetic Reconstruction*. (eds K.A. Joysey and A.E. Friday), pp. 21–74. New York, Academic Press.
- Patterson, C. (1988). The impact of evolutionary theories on systematics. In *Prospects in Systematics* (ed. D.L. Hawksworth), pp. 59–91. Syst. Assoc. Spec. Vol. 36. Oxford, Clarendon Press.
- Pellegrini, M., Marcotte, E.M., Thompson, M.J., Eisenberg, D. and Yeates, T.O. (1999). Assigning protein functions by comparative genome analysis: Protein phylogenetic profiles. *Proc. Natl. Acad. Sci. USA* **96**: 4285–4288.
- Penny, D., Lockhart, P.J., Steel, M.A. and Hendy, M.D. (1994). The role of models in reconstructing evolutionary trees. In *Models in Phylogeny Reconstruction* (eds R.W. Scotland, D.J. Siebert and D.M. Williams), pp. 211–230. Systematics Assoc. Special vol. 52. Oxford, Clarendon Press.
- Penny, D., Hendy, M.D., Lockhart, P.J. and Steel, M.A. (1996). Corrected parsimony, minimum evolution, and Hadamard conjugations. *Syst. Biol.* **45**: 596–606.
- Peterson, K.J. and Eernisse, D.J. (2001). Animal phylogeny and the ancestry of bilaterians: Inferences from morphology and 18S rDNA gene sequences. *Evol. Dev.* **3**: 170–205.
- Phillips, A., Janies, D. and Wheeler, W.C. (2000). Multiple sequence alignment in phylogenetic analysis. *Mol. Phylogenet. Evol.* **16**: 317–330.
- Planet, P.J., DeSalle, R., Siddall, M., Bael, T., Sarkar, I.N. and Stanley, S.E. (2001). Systematic analysis of DNA microarray data: Ordering and interpreting patterns of gene expression. *Genome Res.* **11**: 1149–1155.
- Platnick, N.I. (1979). Philosophy and the transformation of cladistics. *Syst. Zool.* **28**: 537–546.
- Platnick, N.I. and Cameron, H.D. (1977). Cladistic methods in textual, linguistic, and phylogenetic analysis. *Syst. Zool.* **26**: 380–385.
- Platnick, N.I., Griswold, C.E. and Coddington, J.A. (1991). On missing entries in cladistic analysis. *Cladistics* **7**: 337–343.
- Pleijel, F. (1995). On character coding for phylogeny reconstruction. *Cladistics* **11**: 309–315.
- Pol, D. and Siddall, M. (2001). Biases in maximum likelihood and parsimony: A simulation approach to a 10-taxon case. *Cladistics* **17**: 266–281.
- Popper, K. (1957). *The Poverty of Historicism*. London, Routledge and Kegan Paul.
- Popper, K. (1959). *The Logic of Scientific Discovery*. New York, Harper and Row [1968 edition].
- Popper, K. (1962a). Some comments on truth and the growth of knowledge. In *Logic, Methodology and Philosophy of Science* (eds E. Nagel, P. Suppes and A. Tarski), pp. 285–292. Proc 1960 Internatl. Congress. Stanford, CA, Stanford University Press.
- Popper, K. (1962b). *Conjectures and Refutations: The Growth of Scientific Knowledge*. London, Routledge and Kegan Paul.
- Popper, K. (1968). *The Logic of Scientific Discovery*. New York, Harper Torchbooks.
- Popper, K. (1979). *Objective Knowledge. An Evolutionary Approach*. New York, Oxford University Press.
- Popper, K. (1980). Evolution. *New Scientist* **87**: 611.
- Popper, K. (1983). *Realism and the Aim of Science*. London, Routledge.
- Posada, D. and Crandall, K. (1998). MODELTEST: Testing the model of DNA substitution. *Bioinformatics* **14**: 817–818.
- Posada, D. and Crandall, K. (2001a). Selecting models of nucleotide substitution: An application to human immunodeficiency virus 1 (HIV-1). *Mol. Biol. Evol.* **18**: 897–906.
- Pritchard, P.C.H. (1994). Cladistics: The great delusion. *Herpetol. Rev.* **25**: 103–110.
- Posada, D. and Crandall, K. (2001b). Selecting the best-fit model of nucleotide substitution. *Syst. Biol.* **50**: 580–601.
- Prömel, H.J. and Steger, A. (2000). A new approximation algorithm for the Steiner tree problem with performance ratio 5/3. *J. Algorithms* **36**: 89–101.

- Pupko, T., Pe'er, I., Shamir, R. and Graur, D. (2000). A fast algorithm for joint reconstruction of ancestral amino acid sequences. *Mol. Biol. Evol.* **17**: 890–896.
- Pupko, T., Pe'er, I., Hasegawa M., Graur, D. and Friedman, N. (2002). A branch-and-bound algorithm for the inference of ancestral amino-acid sequences when the replacement rate varies among sites: Application to the evolution of five gene families. *Bioinformatics* **18**: 1116–1123.
- Quine, W.V. (1963). On simple theories of a complex world. *Synthese* **15**: 103–106.
- Raff, R.A. (1996). *The Shape of Life: Genes, Development, and the Evolution of Animal Form*. Chicago, IL, University of Chicago Press.
- Rain, J.-C., Selig, L., De Reuse, H., Battaglia, V., Reverdy, C., Simon, S., Lenzen, G., Petel, F., Wojcik, J., Schachter, V. et al. (2001). The protein-protein interaction map of *Helicobacter pylori*. *Nature* **409**: 211–215.
- Rannala B. and Yang, Z. (1996). Probability distribution of molecular evolutionary trees: A new method of phylogenetic inference. *J. Mol. Evol.* **43**: 304–311.
- Reichenbach, H. (1956). *The Direction of Time*. Berkeley, CA, University of California Press.
- Remane, A. (1952). *Die Grundlagen des natürlichen Systems, der vergleichenden Anatomie und der Phylogenetik*. Leipzig, Akademische Verlagsgesellschaft Geest & Portig.
- Resch, A., Xing, Y., Alekseyenko, A., Modrek, B. and Lee, C. (2004). Evidence for a subpopulation of conserved alternative splicing events under selection pressure for protein reading frame conservation. *Nucleic Acids Res.* **32**: 1261–1269.
- Rexová, K., Frynta, D. and Zrzavý, J. (2003). Cladistic analysis of languages: Indo-European classification based on lexicostatistical data. *Cladistics* **19**: 120–127.
- Rice, K.A., Donoghue, M.J. and Olmstead, R.G. (1997). Analyzing large data sets: *rbcL* 500 revisited. *Syst. Biol.* **46**: 554–563.
- Rieppel, O.C. (1988). *Fundamentals of Comparative Biology*. Basel, Birkhäuser Verlag.
- Rieppel, O. (2003). Semaphoronts, cladograms and the roots of total evidence. *Biol. J. Linn. Soc.* **80**: 167–186.
- Rieppel, O. and Kearney, M. (2002). Similarity. *Biol. J. Linn. Soc.* **75**: 59–82.
- Rieseberg, L.H. and Soltis, D.E. (1991). Phylogenetic consequences of cytoplasmic gene flow in plants. *Evol. Trends Plants* **5**: 5–84.
- Rogers, J. (1997). On the consistency of the maximum likelihood estimation of phylogenetic trees from nucleotide sequences. *Syst. Biol.* **46**: 354–357.
- Rogozin, I.B., Wolf, Y.I., Sorokin, A.V., Mirkin, B.G. and Koonin, E.V. (2003). Remarkable interkingdom conservation of intron positions and massive, lineage-specific intron loss and gain in eukaryotic evolution. *Curr. Biol.* **13**: 1512–1517.
- Rokas, A. and Holland, P.W.H. (2000). Rare genomic changes as a tool for phylogenetics. *Trends Ecol. Evol.* **15**: 454–459.
- Romero, I., Fuertes, A., Benito, M.J., Malpica, J. M., Leyva, A. and Paz-Ares, J. (1998). More than 80 R2R3-MYB regulatory genes in the genome of *Arabidopsis thaliana*. *Plant J.* **14**: 273–284.
- Rossnes, R. (2004). *Ancestral Reconstruction of Continuous Characters and its Potential Application to Gene Expression and Alternative Splicing Analysis*. MSc Thesis, University of Bergen, Norway.
- Roth, V.L. (1984). On homology. *Biol. J. Linn. Soc.* **22**: 13–29.
- Roth, V.L. (1988). The biological basis of homology. In *Ontogeny and Systematics* (ed. Humpries, C.J.). New York, Columbia University Press.
- Royall, R. (1997). *Statistical Evidence — a Likelihood Paradigm*. New York, Chapman and Hall.
- Ruse, M. (1971). Narrative explanation and the theory of evolution. *Can. J. Phil.* **1**: 59–74.
- Russell, B. (1948). *Human Knowledge, its Scope and Limits*. London, George Allen and Unwin.
- Rutishauser, R. and Sattler, R. (1989). Complementary and heuristic value of contrasting models in structural biology. III. Case study on shoot-like “leaves” and leaf-like “shoots” in *Utricularia macrorhiza* and *Utricularia purpurea* (Lentibulariaceae). *Botanische Jahrbücher für Systematik* **111**: 121–137.
- Rzhetsky, A. and Nei, M. (1992). A simple method for estimating and testing minimum-evolution trees. *Mol. Biol. Evol.* **9**: 945–967.
- Rzhetsky, A., Ayala, F.J., Hsu, L.C., Chang, C. and Yoshida, A. (1997). Exon/intron structure of aldehyde dehydrogenase genes supports the “introns-late” theory. *Proc. Natl. Acad. Sci. USA* **94**: 6820–6825.
- Saitou, N. and Nei, M. (1987). The neighbor-joining method: A new method for reconstructing phylogenetic trees. *Mol. Biol. Evol.* **4**: 406–425.
- Salem, A.H., Ray, D.A., Xing, J., Callinan, P.A., Myers, J.S., Hedges, D.J., Garber, R.K., Witherspoon, D.J., Jorde, L.B. and Batzer, M.A. (2003). Alu elements and hominid phylogenetics. *Proc. Natl. Acad. Sci. USA* **100**: 12787–12791.
- Salisbury, B.A. (1999). Strongest evidence: Maximum apparent phylogenetic signal as a new cladistic optimality criterion. *Cladistics* **15**: 137–149.
- Salmon, W.C. (1966). *The Foundations of Scientific Inference*. Pittsburgh, PA, University of Pittsburgh Press.
- Sanderson, M. and Hufford, L., eds (1996). *Homoplasy*. San Diego, CA, Academic Press.

- Sanderson, M. and Kim, J. (2000). Parametric phylogenetics? *Syst. Biol.* **49**: 817–829.
- Sanderson, M.J., Purvis, A. and Henze, C. (1998). Phylogenetic supertrees: Assembling the trees of life. *Trends Ecol. Evol.* **13**: 105–109.
- Sankoff, D. (1975). Minimal mutation trees of sequences. *SIAM J. Appl. Math.* **28**: 35–42.
- Sankoff, D. and Blanchette, M. (1998). Multiple genome rearrangement and breakpoint phylogeny. *J. Comp. Biol.* **5**: 555–570.
- Sankoff, D. and Cedergren, R.J. (1983). Simultaneous comparison of three or more sequences related by a tree. In *Time Warps, String Edits, and Macromolecules. The Theory and Practice of Sequence Comparison* (eds D. Sankoff, and J. Kruskal), pp. 253–263. Stanford, CA, CSLI Publications (1999 reprint).
- Sankoff, D. and Nadeau, J.H. (eds) (2000). *Comparative Genomics. Empirical and Analytical Approaches to Gene Order Dynamics, Map Alignment and the Evolution of Gene Families*. Dordrecht, Kluwer Academic Publishers.
- Sankoff, D. and Rousseau, P. (1975). Locating the vertices of a Steiner tree in arbitrary space. *Math. Program.* **9**: 240–246.
- Sankoff, D., Cedergren, R.J. and Lapalme, G. (1976). Frequency of insertion-deletion, transversion, and transition in the evolution of 5S ribosomal RNA. *J. Mol. Evol.* **7**: 133–149.
- Sankoff, D., Morel, C. and Cedergren, R.J. (1973). Evolution of 5S RNA and the non-randomness of base replacement. *Nat. New Biol.* **245**: 232–234.
- Sarkar, I.N., Planet, P.J., Bael, T.E., Stanley, S.E., Siddall, M., DeSalle, R. and Figurski, D.H. (2002). Characteristic attributes in cancer microarrays. *J. Biomed. Inform.* **35**: 111–122.
- Schwikowski, B. and Vingron, M. (1997). The deferred path heuristic for the generalized tree alignment problem. *J. Comput. Biol.* **4**: 415–431.
- Schwikowski, B. and Vingron, M. (2003). Weighted sequence graphs: Boosting iterated dynamic programming using locally suboptimal solutions. *Discr. Appl. Math.* **127**: 95–117.
- Scriven, M. (1959). Explanation and prediction in evolutionary theory. *Science* **130**: 477–482.
- Seitz, V., Ortiz García, S. and Liston, A. (2000). Alternative coding strategies and the inapplicable data coding problem. *Taxon* **49**: 47–54.
- Sellers, P.H. (1974). An algorithm for the distance between two sequences. *J. Comb. Theory* **16**: 253–258.
- Semple, C. and Steel, M. (2002). Tree reconstruction from multi-state characters. *Adv. Appl. Math.* **28**: 169–184.
- Semple, C. and Steel, M. (2003). *Phylogenetics*. Oxford, Oxford University Press.
- Shenkin, P.S., Erman, B. and Mastrandrea, L.D. (1991). Information-theoretical entropy as a measure of sequence variability. *Proteins Struct. Funct. Genet.* **11**: 297–313.
- Sicheritz-Ponten, T. and Andersson, S.G. (2001). A phylogenomic approach to microbial evolution. *Nucleic Acids Res.* **29**: 545–552.
- Siddall, M. (1998). Success of parsimony in the four-taxon case: Long branch repulsion by likelihood in the Farris zone. *Cladistics* **14**: 209–220.
- Siddall, M.E. and Kluge, A.G. (1997). Probabilism and phylogenetic inference. *Cladistics* **13**: 313–336.
- Sikes, D.S. and Lewis, P.O. (2001). *PAUPRat: PAUP Implementation of the Parsimony Ratchet*. Published by the authors (distributed through www.ucalgary.ca/~dsikes/sikes_lab.htm).
- Simmons, M.P. (2004). Independence of alignment and tree search. *Mol. Phylogenet. Evol.* **31**: 874–879.
- Simmons, M.P. and Ochoterena, H. (2000). Gaps as characters in sequence-based phylogenetic analyses. *Syst. Biol.* **49**: 369–381.
- Simon, D. and Larget, B. (1998). *Bayesian Analysis in Molecular Biology and Evolution (BAMBE)*, version 1.01 beta. Pittsburgh, PA, Department of Mathematics and Computer Science, Duquesne University.
- Simpson, G.G. (1964). *This View of Life: The World of an Evolutionist*. New York, Harcourt, Brace and World.
- Sinsheimer, J., Lake, J.A. and Little, R.J.A. (1996). Bayesian hypothesis testing of four-taxon topologies using molecular sequence data. *Biometrics* **52**: 193–210.
- Slatkin, M. and Maddison, W. (1989). A cladistic measure of gene flow inferred from the phylogenies of alleles. *Genetics* **123**: 603–613.
- Slowinski, J.B. (1998). The number of multiple alignments. *Mol. Phyl. Evol.* **10**: 264–266.
- Smith R.L. and Sytsma, K.J. (1990). Evolution of *Populus nigra* (sect. *Aigeiros*): Introgressive hybridization and the chloroplast contribution of *Populus alba* (sect. *Populus*). *Am. J. Bot.* **77**: 1176–1187.
- Smith, T.F., Waterman, M.S. and Fitch, W.M. (1981). Comparative biosequence metrics. *J. Mol. Evol.* **18**: 38–46.
- Smith, V.S., Page, R.D.M. and Johnson, K.P. (2004). Data incongruence and the problem of avian louse phylogeny. *Zool. Scripta* **33**: 239–259.
- Smouse, P.E. and Li, W.-H. (1989). Likelihood analysis of mitochondrial restriction-cleavage patterns for the human-chimpanzee-gorilla trichotomy. *Evolution* **41**: 1162–1176.
- Snel, B., Bork, P. and Huynen, M.A. (1999). Genome phylogeny based on gene content. *Nat. Genet.* **21**: 108–110.
- Snel, B., Bork, P. and Huynen, M.A. (2002). Genomes in flux: The evolution of archaeal and proteobacterial gene content. *Genome Res.* **12**: 17–25.

- Sober, E. (1980). Evolution, population thinking, and essentialism. *Phil. Sci.* **47**: 350–383.
- Sober, E. (1981). The principle of parsimony. *Br. J. Phil. Sci.* **32**: 145–156.
- Sober, E. (1983). Parsimony methods in systematics — philosophical issues. *Annu. Rev. Ecol. Syst.* **14**: 335–357.
- Sober, E. (1985). A likelihood justification for parsimony. *Cladistics* **1**: 209–233.
- Sober, E. (1986). Parsimony and character weighting. *Cladistics* **2**: 28–42.
- Sober, E. (1988a). *Reconstructing the Past: Parsimony, Evolution and Inference*. Cambridge, MA, MIT Press.
- Sober, E. (1988b). The conceptual relationship of cladistic phylogenetics and vicariance biogeography. *Syst. Zool.* **37**: 245–253.
- Sober, E. (1993). *Philosophy of Biology*. San Francisco, CA, Westview Press.
- Sober, E. (1994). Let's razor Ockham's Razor. In *Explanation and its Limits* (ed. D. Knowles), pp. 73–93. *Suppl. Philosophy, Roy. Inst. Phil.* **27**.
- Sober, E. (1996). Parsimony and predictive equivalence. *Erkenntnis* **44**: 167–197.
- Sober, E. (2002). Reconstructing ancestral character states — a likelihood perspective on cladistic parsimony. *The Monist* **85**: 156–176.
- Sober, E. (2003). Parsimony. In *The Philosophy of Science: an Encyclopedia* (eds S. Sarkar and J. Pfeifer). London, Routledge.
- Sober, E. (2004a). The contest between likelihood and parsimony. *Syst. Biol.* **53**: 644–653.
- Sober, E. (in press): Is drift a serious alternative to natural selection as an explanation of complex adaptive traits? In *The Spandrels of San Marco Twenty-Five Years After* (ed. D. Walsh). Oxford, Oxford University Press.
- Sober, E. and Steel, M. (2002). Testing the hypothesis of common ancestry. *J. Theor. Biol.* **218**: 395–408.
- Sokal, R.R. (1986). Phenetic taxonomy: Theory and methods. *Annu. Rev. Ecol. Syst.* **17**: 423–442.
- Sokal, R.R. and Camin, J.H. (1965). The two taxonomies: Areas of agreement and conflict. *Syst. Zool.* **14**: 176–195.
- Soltis, D.E., Soltis, P.S., Chase, M.W., Mort, M.E., Albach, D.C., Zanis, M., Savolainen, V., Hahn, W.H., Hoot, S.B., Fay, M.F. *et al.* (2000). Angiosperm phylogeny inferred from 18S rDNA, *rbcL*, and *atpB* sequences. *Bot. J. Linn. Soc.* **133**: 381–461.
- Sonnhammer, E.L. and Koonin, E.V. (2002). Orthology, paralogy and proposed classification for paralog subtypes. *Trends Genet.* **18**: 619–620.
- Steel, M.A. (1993). Distributions on bicoloured binary trees arising from the principle of parsimony. *Discrete Appl. Math.* **41**: 245–261.
- Steel, M. (2002). Some statistical aspects of the maximum parsimony method. In *Molecular Systematics and Evolution: Theory and Practice* (eds R. De Salle, R. Giribet and W. Wheeler), pp. 125–139. Basel, Birkhäuser Verlag.
- Steel, M. and Penny, D. (2000). Parsimony, likelihood, and the role of models in molecular phylogenetics. *Mol. Biol. Evol.* **17**: 839–850.
- Steel, M. and Penny, D. (2004). Two links between MP and ML under the Poisson model. *Applied Math. Lett.* (in press).
- Steel, M., Penny, D. and Hendy, M. (1993). Parsimony can be consistent! *Syst. Biol.* **42**: 581–587.
- Steel, M., Szekely, L. and Hendy, M. (1994). Reconstructing trees from sequences whose sites evolve at variable rates. *J. Comp. Biol.* **1**: 153–163.
- Steffansson, P. (2004). *Inferring Duplication and Loss Events Using Soft Parsimony*. MSc Thesis, Royal Institute of Technology, Stockholm, Sweden.
- Stevens, P.F. (1984). Homology and phylogeny: Morphology and systematics. *Syst. Bot.* **9**: 395–409.
- Stevens, P.F. (1991). Character states, Morphological variation, and phylogenetic analysis: A review. *Syst. Bot.* **16**: 553–583.
- Storm, C.E. and Sonnhammer, E.L. (2003). Comprehensive analysis of orthologous protein domains using the HOPS database. *Genome Res.* **13**: 2353–2362.
- Strong, E. and Lipscomb, D. (1999). Character coding and inapplicable data. *Cladistics* **15**: 363–371.
- Stuart, J.M., Segal, E., Koller, D. and Kim, S.K. (2003). A gene-coexpression network for global discovery of conserved genetic modules. *Science* **302**: 249–255.
- Suzuki, Y., Glazko, G. and Nei, M. (2002). Overcredibility of molecular phylogenetics obtained by Bayesian phylogenetics. *Proc. Natl. Acad. Sci. USA* **99**: 16138–16143.
- Swofford, D.L. (1984). *PAUP: Phylogenetic Analysis Using Parsimony*. Champaign, IL, Illinois Natural History Survey.
- Swofford, D.L. (1985). *PAUP: Phylogenetic Analysis Using Parsimony*, vers. 2.4. Champaign, IL, Illinois Natural History Survey.
- Swofford, D.L. (1990). *PAUP: Phylogenetic Analysis Using Parsimony*, vers. 3.0. (incl. vers. 3.0s). Champaign, Illinois Natural History Survey.
- Swofford, D.L. (1991). When are phylogeny estimates from molecular and morphological data incongruent? In *Phylogenetic Analysis of DNA Sequences*. (eds M.M. Miyamoto and J. Cracraft), pp. 295–333. New York, Oxford University Press.
- Swofford, D.L. (1993). *PAUP: Phylogenetic Analysis Using Parsimony*, vers. 3.1 (incl. vers. 3.1.1). Champaign, IL, Illinois Natural History Survey.

- Swofford, D.L. (2002). *PAUP*: Phylogenetic Analysis Using Parsimony (*and other methods)*, vers. 4 (incl. vers. 4.0b10). Sunderland, MA, Sinauer Associates.
- Swofford, D.L., Olsen, G.J., Waddell, P.J. and Hillis, D.M. (1996). Phylogenetic inference. In *Molecular Systematics*, 2nd edn (eds D.M. Hillis, C. Moritz and B.K. Marble), pp. 407–514. Sunderland, MA Sinauer Associates.
- Swofford, D., Waddell, P., Huelsenbeck, J., Foster, P., Lewis, P. and Rogers, J. (2001). Bias in phylogenetic estimation and its relevance to the choice between parsimony and likelihoods methods. *Syst. Biol.* **50**: 525–539.
- Tatusov, R.L., Koonin, E.V. and Lipman, D.J. (1997). A genomic perspective on protein families. *Science* **278**: 631–637.
- Tatusov, R.L., Fedorova, N.D., Jackson, J.D., Jacobs, A.R., Kiryutin, B., Koonin, E.V., Krylov, D.M., Mazumder, R., Mekhedov, S.L., Nikolskaya, A.N., *et al.* (2003). The COG database: An updated version includes eukaryotes. *BMC Bioinformatics* **4**: 41.
- Tehler, A., Little, D.P. and Farris, J.S. (2003). The full-length phylogenetic tree from 1 551 ribosomal sequences of chitinous fungi. *Fungi. Mycol. Res.* **107**: 901–916.
- Tellgren, Å., Berglund, A.C., Savolainen, P., Janis, C.M. and Liberles, D.A. (2004). Myostatin rapid sequence evolution in ruminants predates domestication. *Mol. Phylogenet. Evol.* **33**: 782–790.
- Thanaraj, T.A., Stamm, S., Clark, F., Riethoven, J.J., Le Texier, V. and Muilu, J. (2004). ASD: The Alternative Splicing Database. *Nucleic Acids Res.* **32**: D64–D69.
- Thayer, H.S. (1953). *Newton's Philosophy of Nature: Selections from his Writings*. New York, Hafner Publishing Co.
- Thompson, J.D., Higgins, D.G. and Gibson, T.J. (1994). CLUSTAL W: Improving the sensitivity of progressive multiple alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* **22**: 4673–4680.
- Thorne, J.L., Kishino, H. and Felsenstein, J. (1991). An evolutionary model for maximum likelihood alignment of DNA sequences. *J. Mol. Evol.* **33**: 114–124.
- Thorne, J.L., Kishino, H. and Felsenstein, J. (1992). Inching toward reality: An improved likelihood model of sequence evolution. *J. Mol. Evol.* **34**: 3–16.
- Tierney, L. (1994). Markov chains for exploring posterior distributions. *Ann. Stat.* **22**: 1701–1786.
- Tuffley, C. and Steel, M. (1997). Links between maximum likelihood and maximum parsimony under a simple model of site substitution. *Bull. Math. Biol.* **59**: 581–607.
- Uddin, M., Wildman, D.E., Liu, G., Xu, W., Johnson, R.M., Hof, P.R., Kapatos, G., Grossman, L.I. and Goodman, M. (2004). Sister grouping of chimpanzees and humans as revealed by genome-wide phylogenetic analysis of brain gene expression profiles. *Proc. Natl. Acad. Sci. USA.* **101**: 2957–2962.
- Vander Stappen, J., De Laet, J., Gama-López, S., Van Campenhout, S. and Volckaert, G. (2002). Phylogenetic analysis of *Stylosanthes* (Fabaceae) based on the internal transcribed spacer region (ITS) of nuclear ribosomal DNA. *Plant Syst. Evol.* **234**: 27–51.
- Vingron, M. (1999). Sequence alignment and phylogeny construction. In *Mathematical Support for Molecular Biology* (eds M. Farach-Colton, F.S. Roberts, M. Vingron and M. Waterman), pp. 53–64. DIMACS Series in Discrete Mathematics and Theoretical Computer Science, vol. 47. Providence, RI, American Mathematical Society.
- Vrana, P. and Wheeler, W. (1992). Individual organisms as terminal entities: Laying the species problem to rest. *Cladistics* **8**: 67–72.
- Wagner, Jr., W.H. (1952). The fern genus *Diellia*: structure, affinities, and taxonomy. *Univ. Cal. Publ. Bot.* **26**: 1–212, pl. 1–21.
- Wagner, Jr., W.H. (1961). Problems in the classification of ferns. In *Recent Advances in Botany*, vol. 1, pp. 841–844. Toronto, University of Toronto Press.
- Walsh, D. (1979). Occam's razor: A principle of intellectual elegance. *Am. Phil. Q.* **16**: 241–244.
- Wang L. and Jiang, T. (1994). On the complexity of multiple sequence alignment. *J. Comput. Biol.* **1**: 337–348.
- Wang, L., Jiang, T. and Lawler, L. (1996). Approximation algorithms for tree alignment with a given phylogeny. *Algorithmica* **16**: 302–315.
- Wang, L.S., Jansen, R., Moret, B., Raubeson, L. and Warnow, T. (2002). Fast phylogenetic methods for the analysis of genome rearrangement data: An empirical study. *Proceedings of the Pacific Symposium on Biocomputing* (PSB 02), pp. 524–535. Singapore, World Scientific.
- Watanabe, H., Mori, H., Itoh, T. and Gojobori, T. (1997). Genome plasticity as a paradigm of eubacteria evolution. *J. Mol. Evol.* **44**: S57–S64.
- Waterston, R.H., Lindblad-Toh, K., Birney, E., Rogers, J., Abril, J.F., Agarwal, P., Agarwala, R., Ainscough, R., Alexandersson, M., An, P. *et al.* (2002). Initial sequencing and comparative analysis of the mouse genome. *Nature* **420**: 520–562.
- Wheeler, Q.D. (1986). Character weighting and cladistic analysis. *Syst. Zool.* **35**: 102–109.
- Wheeler, Q.D. and Meier, R. (eds) (2000). *Species Concepts and Phylogenetic Theory: a Debate*, pp. 179–184. New York, Columbia University Press.
- Wheeler, W.C. (1994). Sources of ambiguity in nucleic acid sequence alignment. In *Molecular Ecology and Evolution: Approaches and Applications* (eds B. Schierwater,

- B. Streit, G.P. Wagner and R. DeSalle), pp. 323–352. Basel, Birkhäuser Verlag.
- Wheeler, W.C. (1996). Optimization alignment: The end of multiple sequence alignment in phylogenetics? *Cladistics* **12**: 1–9.
- Wheeler, W.C. (1998). Alignment characters, dynamic programming and heuristic solutions. In *Molecular Approaches to Ecology and Evolution* (eds R. DeSalle and B. Schierwater), pp. 243–251. Basel, Birkhäuser Verlag.
- Wheeler, W.C. (1999). Fixed character states and the optimization of molecular sequence data. *Cladistics* **15**: 379–385.
- Wheeler, W.C. (2001a). Homology and the optimization of DNA sequence data. *Cladistics* **17**: S3–S11.
- Wheeler, W. (2001b). Homology and DNA sequence data. In *The Character Concept in Evolutionary Biology* (ed. G.P. Wagner), pp. 303–317. San Diego, Academic Press.
- Wheeler, W.C. (2002). Optimization Alignment: Down, up, error, and improvements. In *Techniques in Molecular Systematics and Evolution* (eds R. DeSalle, G. Giribet and W. Wheeler), pp. 55–69. Basel, Birkhäuser Verlag.
- Wheeler, W.C. (2003a). Implied alignment: A synapomorphy-based multiple-sequence alignment method and its use in cladogram search. *Cladistics* **19**: 261–268.
- Wheeler, W.C. (2003b). Search-based optimization. *Cladistics*, **19**: 348–355.
- Wheeler, W.C. (2003c). Iterative pass optimization of sequence data. *Cladistics* **19**: 254–260.
- Wheeler, W.C. and Gladstein, D.S. (1994). MALIGN: A multiple sequence alignment program. *J. Hered.* **85**: 417–418.
- Wheeler, W.C. and Hayashi, C.Y. (1998). The phylogeny of the extant chelicerate orders. *Cladistics* **14**: 173–192.
- Wheeler, W., Gladstein, D. and De Laet, J. (2003). POY, ver. 3.0.11. Available at <ftp://ftp.amah.org/pub/molecular/poy>.
- Wiley, E.O. (1975). Karl R. Popper, systematics, and classification: A reply to Walter Bock and other evolutionary taxonomists. *Syst. Zool.* **24**: 233–243.
- Wiley, E.O. (1981). *Phylogenetics: the Theory and Practice of Phylogenetic Systematics*. New York, John Wiley and Sons.
- Wilkinson, M. (1995). A comparison of two methods of character construction. *Cladistics* **11**: 297–308.
- Wolf, Y.I., Rogozin, I.B., Grishin, N.V., Tatusov, R.L. and Koonin, E.V. (2001). Genome trees constructed using five different approaches suggest new major bacterial clades. *BMC Evol. Biol.* **1**: 8.
- Wolf, Y.I., Rogozin, I.B., Grishin, N.V. and Koonin, E.V. (2002). Genome trees and the tree of life. *Trends Genet.* **18**: 472–479.
- Wolf, Y.I., Rogozin, I.B. and Koonin, E.V. (2004). Coelomata and not Ecdysozoa: evidence from genome-wide phylogenetic analysis. *Genome Res.* **14**: 29–36.
- Wolfe, K.H. and Sharp, P.M. (1993). Mammalian gene evolution: Nucleotide sequence divergence between mouse and rat. *J. Mol. Evol.* **37**: 441–456.
- Woodger, J.H. (1929). *Biological principles: A critical study*. New York, Harcourt, Brace and Co.
- Wrinch, D. and Jeffreys, H. (1921). On certain fundamental principles of scientific inquiry. *Phil. Mag.* **42**: 369–390.
- Yang, Z. (1994). Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: Approximate methods. *J. Mol. Evol.* **39**: 306–314.
- Yang, Z. (1996). Phylogenetic analysis using parsimony and likelihood methods. *J. Mol. Evol.* **39**: 294–307.
- Yang, Z.H. (1998). Likelihood ratio tests for detecting positive selection and application to primate lysozyme evolution. *Mol. Biol. Evol.* **15**: 568–573.
- Yang, Z.H. and Bielawski, B. (2000). Statistical methods for detecting molecular adaptation. *Trends Ecol. Evol.* **15**: 496–503.
- Yang, Z. and Rannala, B. (1997). Bayesian phylogenetic inference using DNA sequences: A Markov chain Monte Carlo method. *Mol. Biol. Evol.* **14**: 717–724.
- Yang, Z., Goldman, N. and Friday, A. (1995a). Maximum likelihood trees from DNA sequences: A peculiar statistical estimation problem. *Syst. Biol.* **44**: 384–399.
- Yang, Z.H., Kumar, S. and Nei, M. (1995b). A new method of inference of ancestral nucleotide and amino acid sequences. *Genetics* **141**: 1641–1650.
- Yeates, D. (1992). Why remove autapomorphies? *Cladistics* **8**: 387–389.