

Comprehensive Summaries of Uppsala Dissertations
from the Faculty of Science and Technology 937



Bayesian Phylogenetics and the Evolution of Gall Wasps

BY

JOHAN A. A. NYLANDER



ACTA UNIVERSITATIS UPSALIENSIS
UPPSALA 2004

Dissertation presented at Uppsala University to be publicly examined in Lindahlsalen, Evolutionsbiologiskt centrum, Uppsala, Thursday, March 4, 2004 at 13:00 for the degree of Doctor of Philosophy. The examination will be conducted in English.

Abstract

Nylander, J A A. 2004. Bayesian Phylogenetics and the Evolution of Gall Wasps. Acta Universitatis Upsaliensis. *Comprehensive Summaries of Uppsala Dissertations from the Faculty of Science and Technology* 937. 43 pp. Uppsala. ISBN 91-554-5872-6

This thesis concerns the phylogenetic relationships and the evolution of the gall-inducing wasps belonging to the family Cynipidae. Several previous studies have used morphological data to reconstruct the evolution of the family. DNA sequences from several mitochondrial and nuclear genes were obtained and the first molecular, and combined molecular and morphological, analyses of higher-level relationships in the Cynipidae is presented. A Bayesian approach to data analysis is adopted, and models allowing combined analysis of heterogeneous data, such as multiple DNA data sets and morphology, are developed. The performance of these models is evaluated using methods that allow the estimation of posterior model probabilities, thus allowing selection of most probable models for the use in phylogenetics. The use of Bayesian model averaging in phylogenetics, as opposed to model selection, is also discussed.

It is shown that Bayesian MCMC analysis deals efficiently with complex models and that morphology can influence combined-data analyses, despite being outnumbered by DNA data. This emphasizes the utility and potential importance of using morphological data in statistical analyses of phylogeny.

The DNA-based and combined-data analyses of cynipid relationships differ from previous studies in two important respects. First, it was previously believed that there was a monophyletic clade of woody rosid galls but the new results place the non-oak galls in this assemblage (tribes Pediaspidini, Diplolepidini, and Eschatocerini) outside the rest of the Cynipidae. Second, earlier studies have lent strong support to the monophyly of the inquilines (tribe Synergini), gall wasps that develop inside the galls of other species. The new analyses suggest that the inquilines either originated several times independently, or that some inquilines secondarily regained the ability to induce galls. Possible reasons for the incongruence between morphological and DNA data is discussed in terms of heterogeneity in evolutionary rates among lineages, and convergent evolution of morphological characters.

Keywords: Bayesian analysis, Bayes factors, Cynipidae, gall wasps, MCMC, model averaging, model selection, phylogeny, total evidence

Johan A. A. Nylander, Department of Evolutionary Biology, Norbyvägen 18 D, Uppsala University, SE-75236 Uppsala, Sweden

© Johan A. A. Nylander 2004

ISSN 1104-232X

ISBN 91-554-5872-6

urn:nbn:se:uu:diva-3996 (<http://urn.kb.se/resolve?urn=urn:nbn:se:uu:diva-3996>)

In Memory of My Father

List of Papers

This thesis is based on the following papers, which are referred to in the text by their Roman numerals (I–IV):

- I Antonis Rokas, Johan A. A. Nylander, Fredrik Ronquist, and Graham N. Stone. 2002. A maximum-likelihood analysis of eight phylogenetic markers in gall wasps (Hymenoptera: Cynipidae): Implications for insect phylogenetic studies. *Molecular Phylogenetics and Evolution*, 22:206–219.
- II Johan A. A. Nylander, Fredrik Ronquist, John P. Huelsenbeck, and José Luis Nieves-Aldrey. 2004. Bayesian phylogenetic analysis of combined data. *Systematic Biology*, 53:1–21.
- III Johan A. A. Nylander. Phylogenetic inference using Bayesian model averaging: Approximate Methods. *Submitted Manuscript*.
- IV Johan A. A. Nylander, Matthew L. Buffington, Zhiwei Liu, José Luis Nieves-Aldrey, Johan Liljeblad, and Fredrik Ronquist. Molecular phylogeny and evolution of Gall Wasps. *Manuscript*.

Reprint of Paper I was made with the kind permission of the publisher.

In Paper I, AR and JAAN contributed equally with data and analysis. AR completed the text with contributions from JAAN. In Paper II, JAAN did all the laboratory work and analyses, and had major responsibility for the writing of the manuscript, which was completed with assistance from FR. In Paper IV, JAAN provided the majority of the data, conducted the majority of the analyses, and completed the writing of the manuscript with assistance from FR.

Contents

Introduction.....	1
Gall Wasps	1
Methods for Phylogenetic Inference	5
Bayesian Inference of Phylogeny.....	9
Model Selection and Model Averaging.....	10
Objectives	12
Material and Methods	13
Data	13
Analysis.....	13
Results and Discussion	15
Gene Comparisons	15
MCMC Analysis of Combined Data	16
Model Comparisons and Model Averaging	17
The Phylogeny of Gall Wasps.....	18
Evolutionary Implications	20
Convergence or Heterogeneous Evolution in Gall Wasps?	23
Svensk Sammanfattning.....	25
Acknowledgements.....	28
References.....	30

Introduction

"[T]he current trend may be driven primarily by a fascination with the MCMC" (Grant & Kluge, 2003)

Gall Wasps

The most characteristic feature of gall wasps (family Cynipidae) is the capacity of inducing abnormal outgrowths on plants, which are known as galls. When the female wasp lays her eggs, she induces the growth of plant tissue that eventually forms the gall. The gall serves as both protection and nourishment for the developing larva, as it feeds on the gall tissue. Perhaps the most fascinating thing about galls is that they are often abnormal to the plant itself. That is, the gall structure is often wasp-, rather than plant-species specific.

The gall wasps and their closest relatives belong to the Apocrita, the bulk of which are parasitic wasps, and it is very likely that gall wasps originated from insect-parasitic Apocritan forms (Rasnitsyn, 1988; Ronquist et al., 1999). In fact, despite their phytophagous habit, gall wasps can also be viewed as being parasites, only with plant instead of insect hosts. Gall wasps are often parasitized themselves, most often by other parasitic wasps. There are also other species of organisms that use the galls as their home or shelter. This community of organisms — ranging from the plant over the gall wasps, the "guests" feeding on the gall tissue (also called inquilines), and the parasitoids of the gall inducers and inquilines to the parasitoids of the parasitoids (also called hyper-parasitoids) — is a fascinating ecological and evolutionary system well worth studying.

The captivating nature of gall wasps was appreciated early on by scientists such as Alfred Kinsey, who studied their evolution in the beginning of the twentieth century. Kinsey, who later became better known for his studies of human behavior, presented one of the first comprehensive accounts of the phylogeny of gall wasps in 1920 (Kinsey, 1920), well before the strict formalization of quantitative methods for phylogeny reconstruction. More recently, the phylogeny and evolution of gall wasps have received considerable attention (Liljeblad and Ronquist, 1998; Nieves-Aldrey, 2001; Ronquist, 1994; 1995; 1999; Ronquist and Liljeblad, 2001). This work has

led to new insights into the evolutionary history of the group and has, among other things, resulted in a revision of the classification of the gall wasps and their closest insect-parasitic relatives in the superfamily Cynipoidea (Table 1).

Table 1: *Brief overview of the taxonomy, diversity and biology of Cynipoidea, the gall wasps (Cynipidae) and their closest relatives (From Ronquist, 1999; Buffington et al., in preparation).*

Family	Number of genera/species	Distribution	Biology
AUSTROCYNIPIDAE	1/1	Australia	Parasitoids of Lepidoptera larvae
IBALIIDAE	3/19	Holarctic	Parasitoids of Hymenoptera larvae
LIOPTERIDAE	10/170	Widespread, mainly Tropical	Parasitoids of Coleoptera larvae
FIGITIDAE	132/1400	Cosmopolitan	Parasitoids of Diptera, Hymenoptera and Neuroptera larvae
CYNIPIDAE	77/1400	Mainly Holarctic	Phytophagous gall inducers or inquilines
Synergini	7/171	Mainly Holarctic	Phytophagous inquilines in galls of other cynipids
Aylacini	21/156	Holarctic	Gallers on eudicot herbs, one genus also on <i>Smilax</i> vines and <i>Rubus</i> bushes
Diplolepidini	2/63	Holarctic	Gallers on <i>Rosa</i>
Eschatocerini	1/3	South American	Gallers on <i>Acacia</i> and <i>Prosopis</i> (Fabaceae)
Pediaspidini	2/2	Palaearctic	Gallers on <i>Acer</i>
Cynipini	44/974	Mainly Holarctic	Gallers on Fagaceae and Nothofagaceae, mostly on <i>Quercus</i>

More than half of the described species of cynipoid wasps are parasitoids (Table 1). Some of these belong to the so-called "macro cynipoids", a paraphyletic grade comprising the families Austrocynipidae, Ibalidae and Liopteridae. The macro cynipoids attack insect larvae developing in hard substrates (wood, twigs, stems, cones). The rest of the parasitoids belong to the family Figitidae and are micro cynipoids, which are wasps that are generally smaller than macro cynipoids. Figitids attack larvae of Hymenoptera, Diptera and Neuroptera that live in the aphid community, in galls, in decaying organic matter, and in other microhabitats. There are about 1,400 described species of figitids, placed in nine different subfamilies (Ronquist, 1999).

In addition to the Figitidae, the micro cynipoids include the Cynipidae or the "true" gall wasps (Table 1). All cynipids are phytophagous but not all are gall inducers. Within the Cynipidae, there is also a peculiar life mode usually

termed inquilinism. The inquilines have supposedly lost their ability to induce galls but retain the capability of completing galls initiated by other species.

All extant cynipids are placed in a single subfamily, which is divided into six tribes (Ronquist, 1999). Most species belong to the tribe Cynipini, the oak gallers. Other large tribes are the Diplolepidini (gallers of *Rosa*) and the Aylacini, mostly herb gallers. The Pediaspidini (maple gallers) and Eschatocerini (gallers of *Acacia* and *Prosopis*) both include only a few species. Finally, the inquilines are placed in their own tribe, the Synergini. The inquilines attack galls on *Rosa* induced by Diplolepidini, galls on *Quercus* induced by Cynipini, and galls on *Rubus* induced by the Aylacini genus *Diastrophus*. The latter genus is unusual among the Aylacini in attacking a woody host plant but it also includes gallers of rosaceous herbs, such as *Potentilla* and *Fragaria*.

Some authors (e.g., Askew, 1984; Gauld and Bolton, 1988) have regarded the inquilines as a polyphyletic group, with each inquiline species being more closely related to the gall inducer it attacks than to other inquilines. According to this hypothesis, the rose inquilines would be more closely related to rose gallers than to the oak inquilines. Phylogenetic analyses of adult external morphology, however, lend strong support for the idea that the inquilines are monophyletic and had a single origin (Ronquist, 1994; 1995; Liljeblad and Ronquist, 1998) (Figure 1). More specifically, these analyses suggest that the inquilines are most closely related to the Aylacini Rosaceae gallers (*Xestophanes* and *Diastrophus*), initially attacked similar galls, and later radiated to exploit other cynipid host galls (Ronquist, 1994; 1999).

Phylogenetic studies of relationships among other cynipids (Ronquist, 1994; Liljeblad and Ronquist, 1998) based on adult morphological characters show that all other cynipid tribes are monophyletic except the Aylacini, which form an assemblage of basal cynipid lineages (Figure 1). A large group of Aylacini genera mainly associated with Asteraceae and Lamiaceae form a monophyletic group termed the *Isocolus–Neaylax* lineage by Liljeblad and Ronquist (1998). Liljeblad and Ronquist (1998) also found that the tribes Cynipini, Diplolepidini, Pediaspidini, and Eschatocerini together constitute a monophyletic group. They termed this lineage the woody-rosid gallers, since all members attack woody host plants belonging to the rosid clade of eudicots (APG, 2003). Other groups identified in this analysis include the monophyletic *Phanacis–Timaspis* complex and a paraphyletic grade of poppy gallers (*Barbotinia*, *Aylax*, and *Iraella*).

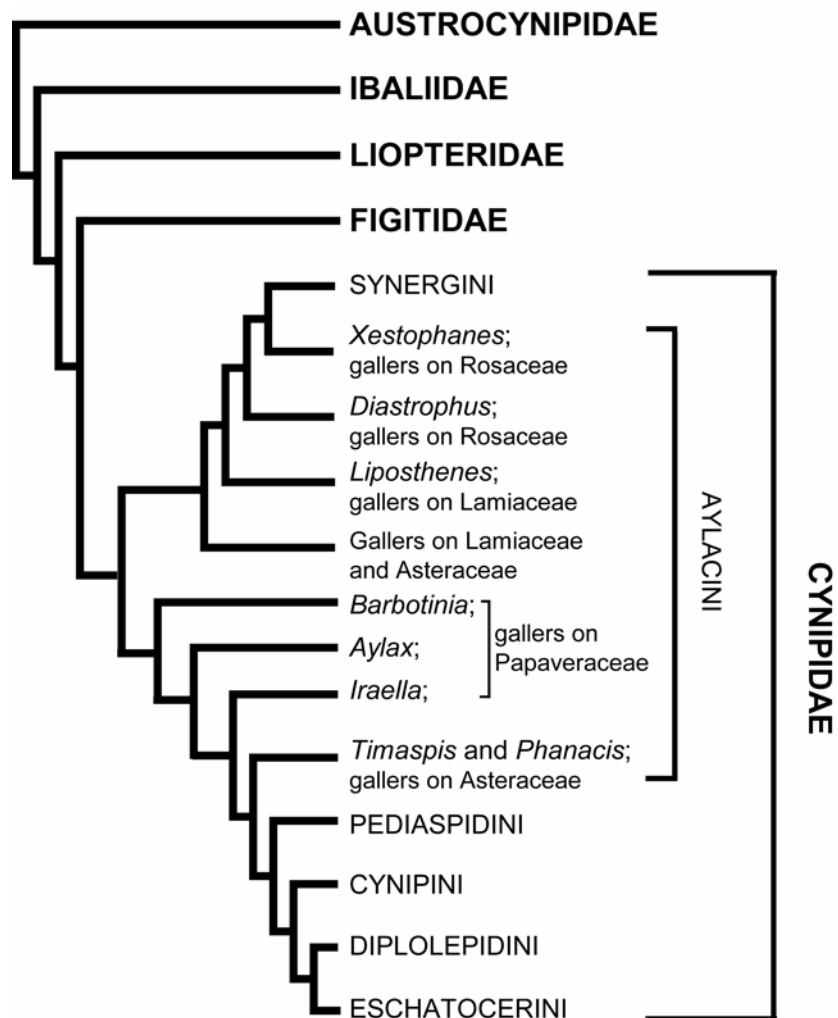


Figure 1. The current view of the phylogeny of the Cynipidae (redrawn from Liljeblad, 2002).

It has long been clear that cynipid gall inducers originated from insect parasitoids. Recent morphology-based analyses of cynipoid phylogeny (Ronquist, 1995; 1999) (Figure 1) indicate that the macrocynipoids form a basal grade leading up to a monophyletic clade including the sister groups Figitidae and Cynipidae. Reconstruction of biological characters on this phylogeny (Ronquist, 1999) suggests that the immediate ancestors of cynipids were parasitoids of Hymenoptera larvae developing inside galls, perhaps gall-inducing Hymenoptera larvae. This is the biology of extant

figitids belonging to the subfamilies Parnipinae and Thrasorinae, which are thought to constitute the earliest figitid lineages.

Until recently, there were two conflicting views on the origin and early evolution of the gall-inducing cynipids. Kinsey (1920) suggested that the first cynipids were herb stem gallers that attacked plants in the Asteraceae and made cryptic galls without visibly affecting the external structure of the plant; he regarded these forms as stem feeders rather than true gallers. Malyshev (1968) considered this unlikely and instead suggested that the first cynipid gallers were associated with oaks or roses and made galls inside seeds. Ronquist and Liljeblad (2001) examined these hypotheses by reconstructing the evolution of life-history traits in Cynipidae based on a morphological phylogeny (Liljeblad and Ronquist, 1998; redrawn in Figure 1), and using non-parametric bootstrapping in combination with parsimony mapping to account for phylogenetic uncertainty. The results showed that the first cynipids with extant descendants induced distinct, single-chambered galls in reproductive organs of herbaceous Papaveraceae, or possibly Lamiaceae. Furthermore, they show that there has been a general trend towards more complex galls in Cynipidae but the herb-stem “feeders” evolved from ancestors inducing distinct galls and their larval chambers are best understood as cryptic galls.

Ronquist and Liljeblad (2001) also analyzed the evolutionary pattern of host-plant usage in cynipids. Gall wasps are conservative in their host-plant choice; indeed, they may be among the most conservative groups of phytophagous insects studied thus far for patterns of host use (Ronquist and Liljeblad, 2001). The evolution of their host-plant preferences is characterized by colonization of pre-existing host-plant lineages rather than by parallel cladogenesis. The few major host shifts, however, have involved remarkably distantly related plant groups. Many shifts have been onto plant species already exploited by other gall wasps, suggesting that interspecific parasitism among cynipids facilitates colonization of novel host plants.

Methods for Phylogenetic Inference

One of the long-standing arguments in systematics has been about which method to use for building phylogenetic trees. Particular focus has been on what kind of optimality criterion to use when searching among all possible candidate trees. Traditionally, simple and intuitive methods, such as using the parsimony method for minimizing steps on trees, have been the method of choice for phylogenetic analysis. Advocates of the parsimony method (MP) have appealed to the philosophical arguments underlying the more general meaning of parsimony, or “Ockham’s razor,” that states that one should prefer simpler explanations (requiring fewer assumptions) over more complex ones. In phylogeny reconstruction, this should correspond to the

feature that MP favours the tree requiring the fewest evolutionary events to explain the observed data and thus is the “simplest” description of the data. Another argument for using MP is that it is claimed to assume as little as possible about the underlying model or mechanism for evolution.

Critics of MP have claimed that the method do make some assumptions about the process underlying the data, and that under some conditions, these assumptions might be unrealistic. These claims appeared when it was shown that MP could be statistically inconsistent (Felsenstein, 1978). That is, MP can, under certain circumstances, lead us to the wrong answer (tree) regardless of the amount of data. These findings spurred the further development of statistical methods for the use in phylogenetics. The method of maximum likelihood (ML), a method that is known to be consistent under a large range of circumstances (e.g., Edwards, 1992; Rogers, 2001), was put forward as an alternative to MP. This was done formally for the use of DNA sequences by Felsenstein (1981), despite being described as a plausible method for phylogenetic inference nearly twenty years earlier (Edwards and Cavalli-Sforza, 1964; Felsenstein, 2003). ML estimation proceeds by assuming a model, and considering the likelihood of a hypothesis (H), given the data (D). The likelihood, $L(H|D)$ is proportional to $P(D|H)$, the conditional probability of observing D given that H is correct (Edwards, 1992). ML inference of phylogeny selects the hypothesis H (the tree) that maximizes the likelihood function for the data D , given a specified model of character evolution (Felsenstein, 1981). The advantage of using ML instead of MP is that while MP only considers the minimal number of changes on a given branch in the tree, ML can "account for" multiple substitutions along that branch. Furthermore, ML allows less probable solutions to potentially influence the results, instead of "placing all bets" on one single answer, as in MP (e.g., Swofford et al., 2001). However, the ML method for phylogenetic inference is computationally demanding and in 1981 systematists lacked the usable methods (software) for applying ML to any but very simple data sets.

Proponents of the MP method have in turn pointed out that despite the fact that the assumption made in statistical methods (ML) are explicit, they are almost always violated, thus consistency cannot be guaranteed with real data (Farris, 1999). For example, the proof of consistency of ML relies upon the fact that the evolutionary model used for inference is the same as the one that generated the data (Rogers, 2001). And since the models applied in ML often are crude approximations to the "true" data-generating model, ML cannot be considered generally consistent. That ML could be inconsistent when the assumed model is strongly violated, has been discussed by a number of authors (e.g., Gaut and Lewis, 1995; Huelsenbeck, 1995; Bruno and Halpern, 1999). This has led to another line of reasoning for the choice of methods based on the notion of robustness. That is, a method is said to robust if it will be relatively unaffected by minor violations of the underlying assumptions (the model). A number of studies have shown (in simulations)

that ML appears to be more robust to violations of the underlying assumptions than e.g., MP, or methods that are based on distance calculations (Huelsenbeck, 1995; Holder, 2001; Sullivan and Swofford, 2001; Swofford et al., 2001). If these observations are extrapolated to real data, it seems that a statistical method would be preferable. Furthermore, the focus on robustness, instead of whether we have, or can ever use the "true" model, lends it support from the theory of model selection based on information theory. That is, scientists do not need to know the "true" model in order to make correct inference or predictions (Burnham and Anderson, 2002). It is still, however, somewhat difficult to conclude that a statistical method will "outperform" other non-statistical methods with real data (Sanderson and Kim, 2000; Steel and Penny, 2000), although "attempting to account for multiple substitutions by using an oversimplified model is a step in the right direction, whereas ignoring them entirely is to accept ignorance." (Swofford et al., 2001:534).

Both ML and MP give point-estimates of the tree (and, for ML, other parameters in the model). To be able to assess the uncertainty in the results, there is a need to use auxiliary criteria. The most commonly used method for assessing confidence is to apply a resampling procedure, such as the bootstrap (Felsenstein, 1985). This means that we need to recalculate the tree (and the parameters) for the resampled data hundreds, or thousands of times. To apply the already time consuming method of ML, which can be orders of magnitude slower than MP, means that it takes too much computational effort to assess confidence in a result. To go around this problem, researchers have taken short-cuts, often by using a simpler model, or doing a less rigorous search for optimal solutions. This can, however, result in the inability of finding an optimal answer, and in the worst case, give biased results (Sanderson and Kim, 2000).

In statistical analysis, an alternative to searching for the single highest point in the "parameter landscape" is to use what is called marginal estimation. Marginal estimation means that we make inference about a particular parameter in our model (such as the tree), while integrating (marginalizing) over the uncertainty in all other parameters. This is, however, potentially even more demanding than finding the maximum of the likelihood, and have not until very recently been possible for real data. Bayesian inference (BI), recently put forward for the use in phylogenetic inference, allows us to do this. BI aims at estimating the posterior probability distribution for a parameter (such as the tree), given the data, a model of character evolution, and prior probabilities on parameters in the model. In brief, the prior is updated in the light of the data to give the posterior, the updated belief in, e.g., a tree. Formally, the posterior probability of a hypothesis H , given the data D , and prior probabilities of the hypothesis $P(H)$ is given by Bayes' rule:

$$P(H|D) = P(D|H)P(H) / \sum_H P(D|H)P(H)$$

The summation (or integration, if H is from a continuous variable) is over the number of possible hypothesis. If the number of hypotheses is large, as it is e.g., when summing over all possible trees for a large number of species, the denominator cannot be calculated analytically. However, the calculation of the posterior probability can be done by approximation, using Monte Carlo simulation, and especially the technique called Markov chain Monte Carlo (Metropolis et al., 1953; reviewed in Gilks et al., 1996). The prior probabilities, $P(H)$, express the uncertainty about the parameters before we observed the data. That is, they represent the subjective or personal probability of the parameters given by the researcher. This is probably the most controversial issue about Bayesian inference since, for example, one person's prior might not be another person's prior (Felsenstein, 2003). The likelihood, $P(D|H)$, is the same as used in ML, and the same models that are used in ML can be used in BI. However, a Bayesian analysis has a number of advantages over ML. The result of a Bayesian analysis reflects in a more intuitive way how researchers might view their results; the probability of the result given the observed data. Moreover, since the result is a set of probability distributions, the researcher receives a direct measure of the uncertainty concerning any parameter in the model. In Bayesian terms, credibility values or credibility intervals can be calculated for any parameter in the model. Finally, the use of MCMC makes Bayesian inference vastly faster than ML. This means that systematists can incorporate complex models into their analysis, without the need to make hazardous trade-offs, and still get robust results with estimates of uncertainty.

The field of systematics, and especially phylogenetic inference, is currently undergoing rapid changes (e.g., Archibald et al., 2003), where new methods are constantly being proposed and added to the phylogenetic toolbox. Most of these methods have fallen into the statistical framework, and for sure, the statistical approach to phylogenetic reconstruction is here to stay (Felsenstein, 2001; Whelan et al., 2001). Among the methods proposed lately, BI has received particular attention, and the power and flexibility of this analytical approach in addressing evolutionary questions has been emphasized (Huelsenbeck et al., 2001; Lewis, 2001; Holder and Lewis, 2003). However, it is important to realize that any method used has its limitations. It is equally important to be aware of under which circumstances a method will fail. Many of these issues can be addressed in the statistical framework, and the issues as how we can interpret support values, such as the credibility values received in BI, or the sensitivity of prior assumptions in BI, etc. will provide challenging areas of research.

Bayesian Inference of Phylogeny

Bayesian inference of phylogeny aims at estimating the posterior probabilities of trees, branch lengths, and other parameters of a character-substitution model (such as the transition:transversion ratio, stationary character-state frequencies etc.). The underlying theory of Bayesian inference of phylogeny has been described by a number of authors (Li, 1996; Mau, 1996; Yang and Rannala, 1997; Larget and Simon, 1999) and was recently reviewed by Huelsenbeck et al. (2001; 2002), Lewis (2001) and Holder and Lewis (2003). Following is a brief description of the method used in this thesis (Paper II; Paper III; Paper IV).

The posterior probability of the i th phylogenetic tree, τ_i , conditioned on the observed data (D) can be obtained using Bayes' formula:

$$P(\tau_i|D) = P(D|\tau_i)P(\tau_i) / \sum_{j=1}^{B(n)} P(D|\tau_j)P(\tau_j),$$

where $P(\tau_i|D)$ is the posterior probability of the tree, given the data and the model, $P(D|\tau_i)$ is the likelihood and $P(\tau_i)$ is the prior probability of the i th tree. The summation given in the formula is over all $B(s)$ trees that are possible for s terminal taxa. In this thesis, we used a flat prior on topology by letting $P(\tau_i) = 1/B(s)$ even though other priors are possible (e.g., Yang and Rannala, 1997). The likelihood function is integrated over all possible values for the branch lengths and parameters in the substitution model. Typically, the posterior probability cannot be calculated analytically. However, the posterior probability of phylogenies can be approximated by sampling trees from the posterior probability distribution. Markov chain Monte Carlo (MCMC) is a method for approximating a complicated surface (the surface in which we are interested is the posterior density) using a simulated Markov chain in which transition probabilities are designed such that the stationary distribution of the chain is the posterior density of interest. In our analyses, we use Metropolis-coupled Markov chain Monte Carlo, (Metropolis et al., 1953; Hastings, 1970; Geyer, 1991) to approximate the posterior probabilities of trees and other parameter values in the models employed. The algorithm has been implemented in a software package, MrBayes (Huelsenbeck and Ronquist, 2001; Ronquist and Huelsenbeck, 2003) which was used for all calculations.

Theory predicts (Tierney, 1994) that a Markov chain that is properly set up and run for an infinite number of generations will produce a valid sample from the target distribution. However, in every particular analysis, these criteria are hard to fulfill. In samples drawn by a proper MCMC, the frequency of occurrence is in direct proportion to the posterior probability of a topology or a parameter value. Obviously, it is of utmost importance that the MCMC has reached its target distribution, and that it gives a correct

representation of it, in order to interpret the frequencies as posterior probabilities. Hence, criteria for assessing that the MCMC has reached its target distribution are important (Huelsenbeck et al., 2002; Paper II). Extensive post-run analyses of MCMC are seldom seen in phylogenetics (for an exception, see Drummond and Rambaut, 2003), and more work in these area is much needed.

The result from a MCMC analysis is a set of probability distributions. It is then a question of how this information should be summarized. A commonly used method is to present the mode of a distribution, and perhaps present a credibility region around it. A problem occurs if the distribution is not unimodal. That is, the representation of the variation in probability over the parameter space with a single value might be inappropriate. Optimally, the whole distribution would be given as the result. A commonly used method for summarizing the relationships of taxa is to present a majority-rule consensus tree, and to plot the frequency of occurrence for the branches to represent the posterior probability of clades. Whether this is a valid procedure, that is, if those frequencies are really reflecting the posterior probability of groups is a matter of debate (e.g., Suzuki et al., 2002; Wilcox et al., 2002; Erixon et al., 2003; Ronquist and Huelsenbeck, in press), and will be one of the most important areas of research in the future of Bayesian phylogenetics.

Model Selection and Model Averaging

With the advent of Monte Carlo integration in phylogenetics and especially the method of Bayesian inference using MCMC, analyses using parameter-rich models are now made more feasible. It is now relative easy to implement complex models (such as combining morphology and nucleotides; Paper II; Paper IV), and to run in reasonable time. There is, however, a potential danger in using an increasing number of parameters in models. Increasing complexity comes with a cost of increased error variance of the parameter estimates and eventually leads to what is known as overparameterization or overfitting (Burnham and Anderson, 2002). Potentially, this is problematical for phylogenetic inference since the tree topology itself is a parameter in model-based methods of inference. Hence, there is a trade-off between bias and efficiency when it comes to choosing a model for inference (e.g., Zucchini, 2000; Burnham and Anderson, 2002).

Several methods for comparing and selecting models for phylogenetic inference have recently been discussed. Most common is the use of likelihood ratio tests (LRT) (e.g., Felsenstein, 1981; Huelsenbeck and Crandall, 1997; Posada and Crandall, 2001), where the model is allowed to be made more complex only if the addition of a parameter provides a significant increase in the likelihood. Other authors (e.g., Kishino and

Hasegawa, 1989; Buckley et al., 2002; Pupko et al., 2002) have taken the information-theory based approach using the Akaike information criteria (AIC; Akaike, 1973), where the likelihood score under a model is, in effect, penalized for its model complexity, and the model with the lowest AIC is chosen. Bayesian approaches have been the use of posterior predictive P values, where the accuracy of a model is investigated using simulations (Bollback, 2002), or the use of the Bayesian information criterion (BIC; Schwartz, 1978). BIC resembles AIC insofar as it penalizes the model for its model complexity, and was used by e.g. Posada and Crandall (2001) for comparing substitution models. The use of Bayes factors (Kass and Raftery, 1995) is yet another Bayesian alternative for comparing models. Bayes factors measure the strength of the support for one model over another by taking the ratio of the marginal likelihoods and can be applied for comparing non-nested models (unlike e.g. LRT). The marginal likelihood (also called the predictive, model or integrated likelihood) is simply the likelihood of the data under a given model. In Bayesian inference of phylogeny, the marginal likelihood is given by the denominator of Bayes theorem for calculating posterior probabilities of trees, and can be approximated using the output of an MCMC (Newton and Raftery, 1994; Paper II; Paper III)

The use of a model selection criterion can also help us rank models (except when using the LRT) and find out how much better one model is compared to the rest by calculating "Akaike weights" for individual models (Burnham and Anderson, 2002). But what if there are one or several models that are almost equally good as the one chosen as best under our criterion? Furthermore, depending on the "luck of the draw" the data might point to a suboptimal model that would make inaccurate predictions or inferences from the data. This indicates the presence of model selection uncertainty. How do we take model selection uncertainty into account in the inference? One solution to this is to abandon the idea that we need to select a single model for inference. Instead, we can use all models, and let each model in a set of competing models contribute in proportion to its posterior probability. This approach is called model averaging, and the Bayesian framework allows us to calculate the posterior probabilities of models, and to base inference on a set of models weighted by their (posterior) probability (Wasserman, 2000; Paper III).

Objectives

This thesis concerns the phylogenetic relationships and the evolution of the gall-inducing wasps belonging to the family Cynipidae, and their closest relatives in the superfamily Cynipoidea. I present a first assessment of the higher-level relationships of the gall wasps based on molecular, as well as molecular and morphological data. It is also the first study to do this within a Bayesian framework using new methods for combining data. The main question, and the reason this project was initiated, was:

What do molecular data tell us about the higher-level phylogeny and evolution of gall wasps?

In order to lay out the ground for an answer to this question, we designed a pilot study (Paper I), in which we were interested in:

Which DNA-sequence regions (genes) are useful for inferring phylogenetic relationships in cynipoid wasps?

By utilizing the results in the first paper, we expanded our taxon sample and set out to perform a larger analysis (Paper IV). There was also the question on how the results from molecular and morphological would differ, that is,

Are molecular and morphological data sampled from cynipoids congruent?

For answering this question, we wanted to develop a framework, or a method that could incorporate the different data sets (multiple DNA data sets and morphology) in a single, parametric analysis. Thus, the question (Paper II) was

How can we incorporate models that allow data heterogeneity into Bayesian analysis?

Related to the issue of the performance of model-based inference of combined data was the performance of models for character evolution, and how we could compare and choose between them. Especially, we were interested in (Paper II; Paper III):

How do Bayesian methods for selecting and comparing models perform in phylogenetic analysis?

Material and Methods

Data

Gall wasp specimens for the use as representatives in phylogenetic analyses were collected using nets, Malaise traps or more commonly, by collecting galls and rearing out the adults. In a few cases, galls were opened and the larvae or pupae were collected and preserved in 96% ethanol.

From the exemplar specimens, DNA was extracted and gene fragments were amplified using PCR techniques. We used a range of nuclear and mitochondrial genes, both ribosomal and protein-coding, that previously been used in insect systematics (for a review, see Caterino et al., 2000). The gene fragments were sequenced according to the protocols in Paper I. These were the 18S rDNA (18S), 28S rDNA (28S), 5.8S rDNA (5.8S), ITS1 and ITS2, elongation factor 1- α F1 (EF1 α), long-wavelength opsin (LWRh), cytochrome b (Cytb), and cytochrome oxidase I (COI). For the larger analyses in Paper II and Paper IV, as well as for the comparisons made in Paper III, only some of these gene fragments were used.

The morphological characters used in the combined (also called simultaneous or total-evidence) analyses (Paper II; Paper IV) were taken from Liljeblad and Ronquist (1998) and Ronquist (1999), with some corrections and additions. In Paper IV, life-history ("biological") characters were taken from Ronquist (1999) and Ronquist and Liljeblad (2001), complemented with additional information for some taxa.

Analysis

Phylogenetic trees from DNA sequences (Paper I; Paper II; Paper III; Paper IV), morphological characters (Paper II; Paper IV), and both character sources combined (Paper II; Paper IV) were inferred by parsimony (Paper I; Paper II; Paper III; Paper IV) and maximum likelihood (Paper I) using the program PAUP* (Swofford, 2002). Bayesian analysis using MCMC (Paper II; Paper III; Paper IV) was carried out using the program MrBayes (Huelsenbeck and Ronquist, 2001; Ronquist and Huelsenbeck, 2003).

A range of different criteria was used for comparing and choosing models for character analysis. These included LRT (Paper I), AIC and BIC (Paper II; Paper III), and Bayes factors based on marginal likelihood estimation

(Paper II; Paper III; Paper IV). Model averaging was accomplished by estimating the posterior model probabilities using marginal likelihood estimators and Bayes factors, and then sampling trees from the output of MCMC analyses run under individual models in proportion to their model probability (Paper III).

In Paper IV, the evolution of cynipid life-history characters was reconstructed using Bayesian methods. Character evolution was assumed to follow a discrete-state Markov model, and inferences about ancestral states were drawn using Bayesian MCMC techniques while accounting for the uncertainty in topology, as well as the uncertainty in parameters in the model of character evolution (Schultz and Churchill, 1999; Huelsenbeck et al., 2000; Huelsenbeck and Bollback, 2001; Huelsenbeck and Imennov, 2002).

Results and Discussion

Gene Comparisons

In Paper I, we assessed the utility of eight gene fragments (5.8S, 18S, 28S, ITS 1 and 2, LWRh, EF1 α , Cytb, and COI) in reconstructing phylogenetic relationships at various levels of divergence in gall wasps, using a set of eight exemplar taxa. The result of phylogenetic analyses of the individual loci using ML is shown in Figure 2. Likelihood ratio testing was used to find the best fitting evolutionary model of each of the markers. The likelihood model best explaining the data was, for most loci, parameter-rich, with strong A-T bias for mitochondrial loci and strong rate heterogeneity for the majority of loci. Our data suggest that 28S, EF1 α , and LWRh may be potentially useful markers for the resolution of cynipid and other insect within-family-level divergences (ca. 50–100 mya old), whereas mitochondrial loci and ITS regions might be most useful for lower-level phylogenetics. In contrast, the 18S is likely to be useful for the resolution of above-family-level relationships. For further study of cynipid relationships, we focused on 28S, EF1 α , LWRh, and COI. This set of markers included both nuclear and mitochondrial genes, and the combined signal should enable resolution of old to intermediate splits in the Cynipidae phylogeny.

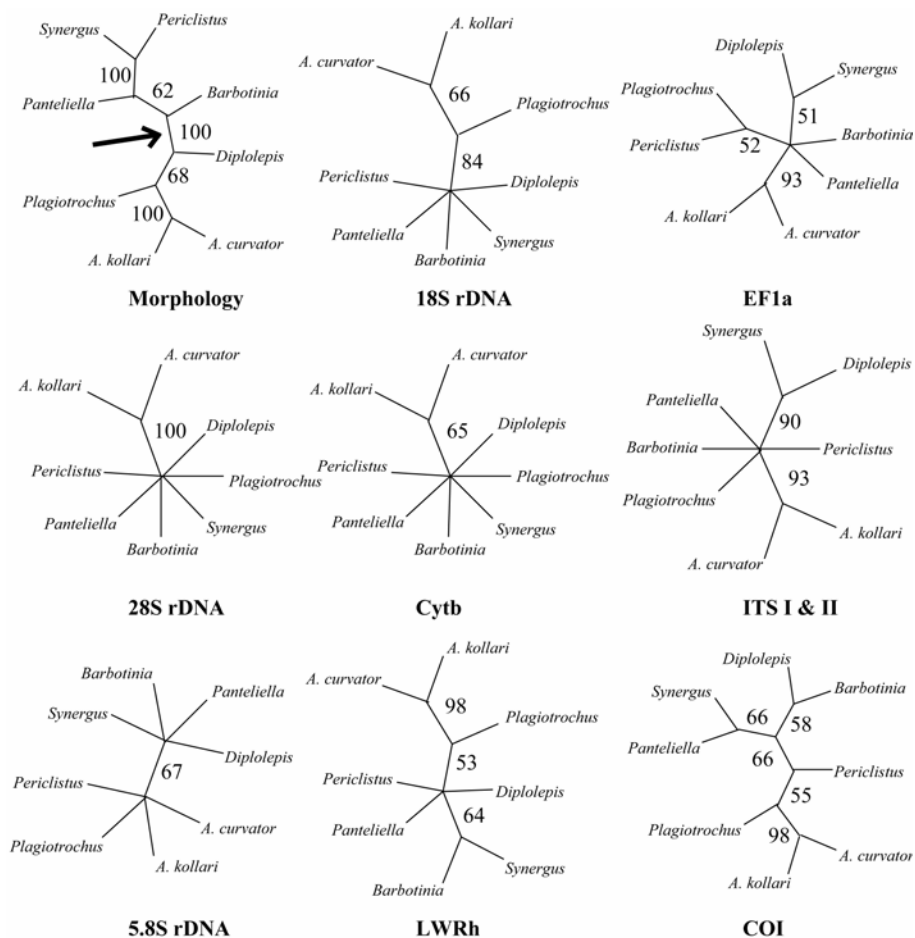


Figure 2 Figure showing the phylogenetic signal in maximum-likelihood analyses of different gene-regions for eight gall-wasp taxa (Paper I). Numbers on branches indicate the bootstrap frequency of occurrence. The morphological tree with parsimony bootstrap frequencies is given as a comparison. Trees are unrooted and the preferred placement of the root is indicated with an arrow in the morphological tree.

MCMC Analysis of Combined Data

In Paper II, we develop a Bayesian MCMC approach to the analysis of combined datasets and explore its utility in inferring relationships among gall wasps based on data from morphology and four genes (COI, LWRh, EF1 α , and 28S). To our knowledge, this is the first statistical phylogenetic analysis of combined morphological and molecular data. Examined models range in complexity from those recognizing only a morphological and a

molecular partition to those having complex substitution models with independent parameters for each gene. We find that Bayesian MCMC analysis deals efficiently with complex models: convergence occurs faster and more predictably for complex than for less complex models, mixing is adequate for all parameters even under very complex models, and the parameter update cycle is virtually unaffected by model partitioning across sites. Morphology contributes only 5 % of the characters in the dataset but nevertheless influences the combined-data tree, supporting the utility of morphological data in multi-gene analyses.

Model Comparisons and Model Averaging

The examined models in Paper II range in complexity from those recognizing only a single partition to those having complex substitution models with independent parameters for each gene. Bayes factors show that process heterogeneity across data partitions is a significant model component, although not as important as among-site rate variation. More complex evolutionary models are associated with more topological uncertainty and less conflict between morphology and molecules. Bayes factors sometimes favour simpler models over considerably more parameter-rich ones but the best model overall is also the most complex one and Bayes factors do not support exclusion of apparently weak parameters from this model. Thus, Bayes factors appear to be useful for selecting among complex models but it is still an open question whether their use strikes a reasonable balance between model complexity and error in parameter estimates.

In Paper III, Bayesian model averaging to phylogenetic inference was introduced. In particular, ways of approximating the posterior probability of models using Akaike weights were compared to those derived by Bayes factors. In this case, it was shown that despite the presence of model selection uncertainty in the data examined, the use of model averaging had only minor effects on phylogeny estimation compared to inference based on a single best model. This was probably due to the fact that the models receiving high posterior model probability were very similar and therefore produced similar phylogenetic results. Furthermore, Akaike weights based on the BIC or the AIC were found to be unreliable estimators of posterior model probabilities, using estimates based on MCMC output as the standard of reference. However, they may still be useful indicators of the models that are likely to have high posterior probability among a set of candidate models considered for model averaging. These models can then be selected for more detailed analysis, an approach known as "Occam's window" (Madigan and Raftery, 1994). Thus, model averaging is a natural extension of the Bayesian framework of phylogenetic analysis and it is demonstrated that it can be accomplished by using methods that estimate marginal likelihoods (Paper

III) and posterior probabilities. Furthermore, this can be done using existing software and for practically any set of models, including those which incorporate partitioning of the data (Paper II).

The Phylogeny of Gall Wasps

In Paper IV, we extended the analysis pioneered in Paper II to include 89 taxa of cynipids and related parasitoids. We sequenced parts of three genes (28S, COI, and EF1 α) for 89 species of cynipids and related parasitoids. The sample included 70 cynipid species, representing the majority of the described cynipid genera in all tribes except the Cynipini, which are generally considered a monophyletic clade and were only represented by a small number of the about 45 described genera. The outgroup sample included representatives from all families in the Cynipoidea except the Austrocynipidae, an Australian endemic only known from three specimens. The species-rich and diverse family Figitidae, considered to be the sister-group of the Cynipidae (Ronquist, 1999), was represented by seven of the nine subfamilies. The molecular data were analyzed separately and in combination with previously published morphological and life-history data (Paper II; Liljeblad and Ronquist, 1998; Ronquist, 1999; Ronquist and Liljeblad, 2001).

The analyses in Paper II and Paper IV are largely congruent with respect to cynipid relationships (the result from Paper IV is given in Figure 3). The DNA data support several of the previous conclusions based on morphological data (Ronquist, 1994; Liljeblad and Ronquist, 1998) but they also conflict with many of them. The molecular data support the monophyly of the oak gallers (Cynipini) and the rose gallers (Diplolepidini) but suggest that the herb gallers (Aylacini) are not monophyletic. These results all agree with previous analyses (cf. Figure 1). However, there are two major differences between the morphological and molecular data. The first concerns the inquilines, the other the woody-rosid gallers.

Adult morphology lends strong support to the monophyly of the inquilines (Synergini) but the DNA analyses (Paper II; Paper IV) split the inquilines into three separate groups: (1) the *Synergus* complex (including *Rhoophilus*) of (largely) oak inquilines; (2) the genus *Ceroptres* of oak inquilines; and (3) the inquilines in Rosaceae galls (*Periclistus* and *Synophromorpha*), which group with the Rosaceae gallers *Diastrophus* and *Xestophanes* nested among them. The molecular data are somewhat inconclusive regarding the relationships among these three groups but indicate that they might have separate origins, and appear as a grade close to the base of the Cynipidae tree. When the molecular data are combined with morphological data, the three inquiline groups end up in a single monophyletic clade, still with the Rosaceae gallers *Xestophanes* and

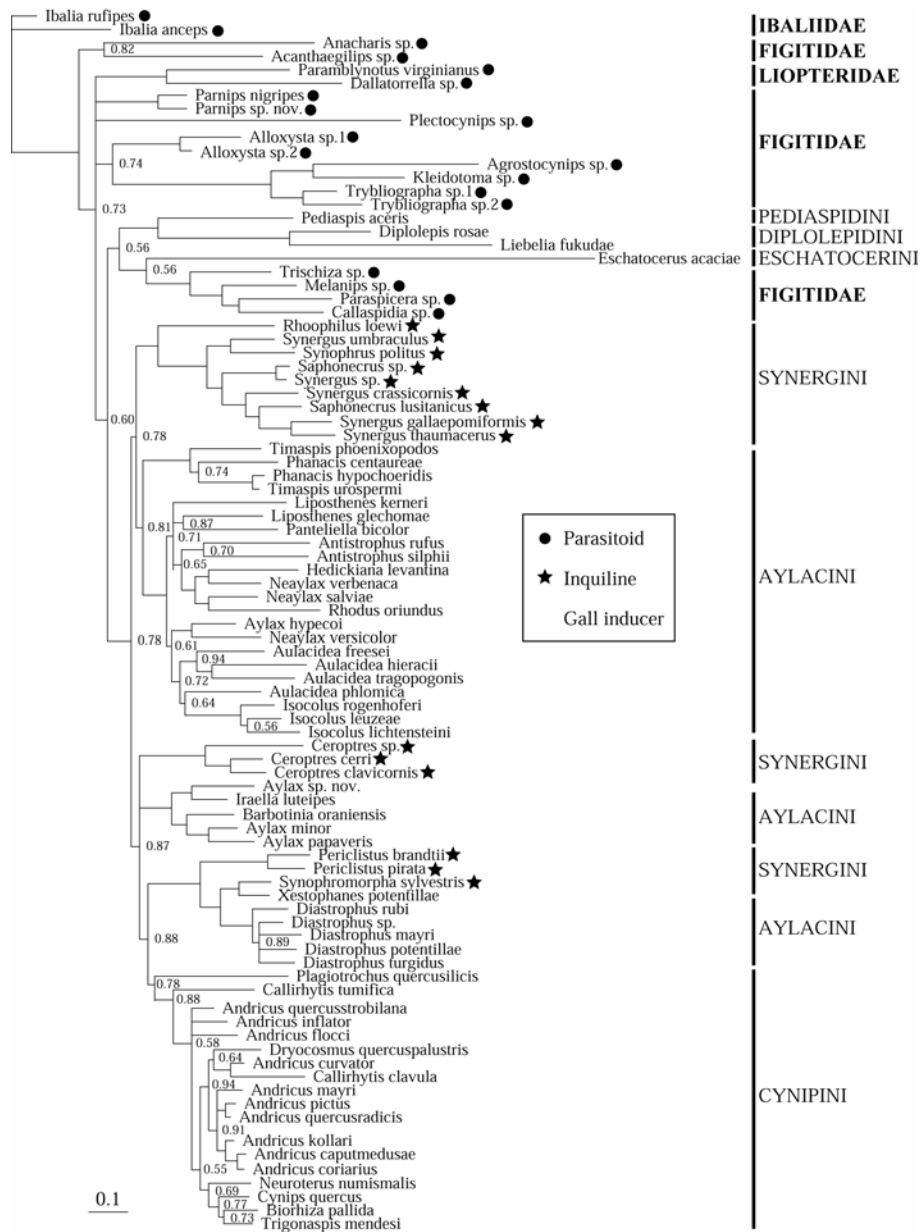


Figure 3 Majority-rule consensus tree based on Bayesian MCMC analysis of DNA data (Paper IV). Life modes (parasitoid, gall inducer, or inquiline) are given along with the family-level classification. Tribes listed belong to Cynipidae. Numbers on branches indicate clade-credibility values (only values between 0.50 and 0.95 are given), and the scale indicate branch-length expresses as the expected number of substitutions.

Diastrophus nested within this clade. The latter result was unexpected, even though morphological data suggest that *Xestophanes* and *Diastrophus* are the gall inducers that are most closely related to the inquilines.

The second conflict concerns the woody-rosid gallers (Diplolepidini, Eschatocerini, Pediaspidini, Cynipini), which all induce galls on woody members of the rosid clade of eudicots. Morphological analyses indicate that they form a monophyletic clade with a single origin (Figure 1) (Ronquist, 1994; Liljeblad and Ronquist, 1998). However, the DNA sequences of the woody-rosid gallers that are not associated with oaks (Diplolepidini, Eschatocerini, and Pediaspidini) are distant from each other and from the sequences of the oak gallers (Cynipini) and other cynipids. The molecular analyses place the non-oak woody-rosid gallers basal to the rest of the Cynipidae. In this case, the molecular signal is strong enough that it remains unaffected by the addition of morphological data (Paper II; Paper IV).

At lower levels, there is much congruence between morphology and molecules. For instance, the DNA data and combined analyses support a clade of Asteraceae and Lamiaceae gallers recognized previously as the *Isocolus-Neaylax* clade (Liljeblad and Ronquist, 1998), albeit with a slightly different circumscription (including rather than excluding the genus *Liposthenes*). The molecular and combined-data circumscription of this clade is supported by a recently described larval character (Nieves-Aldrey, Vårdal, and Ronquist, submitted; Vårdal, 2004). The DNA and combined data further support the monophyly of the *Phanacis-Timaspis* complex of Asteraceae gallers, in agreement with previous studies, and the monophyly of a clade of poppy gallers, previously thought to be a paraphyletic grade.

The phylogenetic results presented in this thesis (Paper II; Paper IV) and earlier (Ronquist, 1994; Liljeblad and Ronquist, 1998) will necessitate a revision of the current classification of the Cynipidae. However, this is deferred until additional data are available to corroborate some of the unexpected groupings emerging in the analyses presented here (Paper II, Paper IV), and to more robustly resolve the basal splits in the Cynipidae tree.

Evolutionary Implications

Evolutionary reconstructions based on morphological phylogenetic analyses (Ronquist, 1994; 1999; Ronquist and Liljeblad, 2001) suggest that the first cynipids were gall-inducers rather than inquilines, and induced single-chambered, distinct, integral swellings in fruits or other reproductive structures of herbs belonging to the family Papaveraceae or possibly Lamiaceae. They apparently originated from internal parasitoids of Hymenoptera larvae developing inside galls (Ronquist, 1999). The geographical center of origin was probably in the Palearctic (Ronquist and Liljeblad, 2001; Table 2).

Table 2: Table summarizing the inferred ancestral states of life-history characters in the Cynipidae based on morphology (Ronquist and Liljeblad, 2001), DNA (Paper IV), and both data combined (Paper IV). The character states are inferred using Bayesian methods except for the morphology column (Morph.), where numbers are from bootstrapped parsimony reconstructions (Ronquist and Liljeblad, 2001). The highest posterior probability for each column and character state is indicated in bold face.

Character	Character state	Posterior prob. for ancestral char. state in Cynipidae		
		Morphology	DNA	Morphology + DNA
Host family	Fagaceae	0.04	0.81	0.71
	Rosaceae	<0.01	0.08	<0.01
	Anacardiaceae	<0.01	<0.01	<0.01
	Lamiaceae	0.26	<0.01	0.02
	Asteraceae	0.04	0.03	0.15
	Papaveraceae	0.66	0.01	0.11
	Fabaceae	<0.01	<0.01	<0.01
	Sapindaceae	<0.01	0.04	<0.01
Host plant form	woody	0.02	0.97	0.54
	herbaceous	0.98	0.03	0.46
Gall chambers	single	0.84	0.77	0.70
	many	0.16	0.23	0.30
Gall position	fruit, seed, flower	0.88	0.43	0.62
	bud	<0.01	0.09	0.01
	leaf	<0.01	0.11	0.05
	stem, twig, runner	0.12	0.31	0.32
	root	<0.01	0.06	<0.01
Gall structure	cryptic	<0.01	0.02	<0.01
	swelling	>0.99	0.96	>0.99
	complex	<0.01	0.02	<0.01
Gall attachment	integral	>0.99	>0.99	>0.99
	semi detachable	<0.01	<0.01	<0.01
	detachable	<0.01	<0.01	<0.01
Distribution	Palaearctic	>0.99	>0.99	>0.99
	Nearctic	<0.01	<0.01	<0.01
	Neotropic	<0.01	<0.01	<0.01
	Australia	<0.01	<0.01	<0.01
	South Africa	<0.01	<0.01	<0.01

With respect to their host-plant relationships, gall wasps are among the most specialized (host-plant specific) and conservative groups of phytophagous insects we know. There are few major shifts between host plants;

nevertheless, there have been cases of independent radiation onto the same set of distantly related host plant species (Ronquist and Liljeblad, 2001). Morphological data have been somewhat ambiguous concerning the origin of alternating generations in cynipids. Some analyses place the two clades with this trait (Pediastidini and Cynipini) as sister groups, suggesting that this complex trait could have a single origin (Liljeblad and Ronquist, 1998). However, other analyses (Liljeblad, 2002) indicate that these clades are more distant, making it more likely that they evolved alternating generations independently.

The inquilines are strongly supported as a monophyletic group by morphological analyses, and they appear most closely related to the Aylacini Rosaceae galls *Diastrophus* and *Xestophanes* (Ronquist, 1994; Liljeblad and Ronquist, 1998). The result from Ronquist and Liljeblad (2001) suggest that the inquilines had a single origin from Aylacini Rosaceae galls, originally attacked similar galls, and later radiated to exploit Cynipini and Diplolepidini galls on Fagaceae and Rosaceae.

The molecular and combined analyses presented here partly confirm these patterns and partly suggest other possibilities. The patterns of conserved host-plant affiliations, rare shifts, and convergent radiation remain supported. The likely ancestral states of cynipids (Table 2) indicate, as before, that the cynipids originated in the Palearctic from Hymenoptera parasitoids, and that the first forms were gall-inducers rather than inquilines. The first galls were likely to have been single-chambered, distinct and integral swellings of fruits or other reproductive structures of plants, again in agreement with previous results. However, the new analyses suggest that the original host plant might well have been a woody plant instead of a herb, most likely an oak or another woody plant in the family Fagaceae (Table 2). The new analyses also increase the probability that the ancestral gall was multi-chambered and induced in stems of Asteraceae, as proposed by Kinsey (1920). However, the probability that the ancestral gall was cryptic remains small (Table 2), contradicting Kinsey's views.

The molecular and combined analyses clearly indicate that alternating generations evolved independently in Pediastidini and Cynipini because of the distant position of these groups in the phylogeny (Figure 3).

Regarding the inquilines, the molecular and combined analyses suggest slightly different scenarios, even though none of them support the monophyly of the inquilines. The molecular results place the three clades of inquilines (one with Rosaceae galls included) separate from each other at the base of the tree. Each of these three clades receives high support, suggesting that there have been at most three different origins of inquilines (or four if *Periclistus* and *Synophromorpha* in the mixed inquiline-gall inducer clade evolved inquilinism independently of each other). It is worth noting that none of the inquiline clades groups strongly with their host clade (Cynipini for *Ceroptres* and most of the *Synergus* complex, Diplolepidini for

Periclistus and *Diastrophus* for *Synophromorpha*), as postulated by the polyphyletic origin hypothesis.

The combined result group all inquilines in a single clade, with the Rosaceae gallers *Diastrophus* and *Xestophanes* deeply nested within. This tree suggests that the inquilines might have had a single origin, after all, but that some lineages apparently reverted to inducing their own galls. An observation that seems to lend some support to this idea is that there is fighting at oviposition sites in *Diastrophus* (Ronquist, 1994; 1999 and references cited therein). This could possibly be associated with intra-specific inquilinism that could be a remnant of an ancestral inquiline life mode.

Convergence or Heterogeneous Evolution in Gall Wasps?

Perhaps the most intriguing aspect of the results presented in this thesis is the strong conflict between morphology and molecules concerning the relationships among inquilines and woody rosid gallers. What is the cause of this pattern? An answer that suggests itself is that it is due to morphological convergence among inquilines and woody-rosid gallers due to their shared life history. Ronquist (1994) developed a technique to examine this possibility and applied it to the origin of the cynipid inquilines. The results suggested it was unlikely that the morphological support for inquiline monophyly was due entirely to convergences in conflict with true relationships. A possibility that is available now is to develop probabilistic models of morphological convergence and examine them in a Bayesian framework.

It is also possible that the conflict is entirely or partly due to imperfections in current models of molecular evolution used in phylogenetic inference. For instance, standard models assume that the evolutionary process is homogeneous across the tree. If the process actually differs in different parts of the tree, it is possible that unrelated branches may erroneously “attract” each other just because they have converged with respect to how their genes evolve, much like long branches attract each other in parsimony analysis. In Paper IV, we examine one model that allows variation across the tree, the so-called covarion-like model (Tuffley and Steel, 1998; Huelsenbeck, 2002). The results of these analyses, however, were similar to those obtained with standard models. Of course, the covarion-like model only represents one way in which the evolutionary process itself can evolve over the tree.

In conclusion, it is still not entirely clear why there is a conflict between the molecular and morphological data concerning the woody-rosid gallers

and inquilines. However, this should prove a worthwhile subject to explore in future research using more data and more sophisticated models and statistical techniques.

Svensk Sammanfattning

Den här avhandlingen handlar om släktskapen mellan olika arter och släkten av gallsteklar inom familjen Cynipidae. Avhandlingen baserar sig på fyra delarbeten (refererade som manuskript I–IV i avhandlingen), som dels beskriver utvecklandet av generella metoder för att analysera släktskap bland organismer i allmänhet, och dels analyserar släktskapen bland gallsteklar i synnerhet.

Steklarna i familjen Cynipidae är mest kända för sin förmåga att bilda galler på växter. De orsakar bland annat det vi kallar galläpplen på ekar och sömntornsgaller på rosenbuskar. När stekelhonan lägger sina ägg i växten påverkar hon samtidigt växtcellerna så att de tillväxer till en skyddande kammare (gall) runt ägget. Gallen erbjuder både skydd och föda för stekeln under dess utveckling då larven äter av innanmätet. Det kanske mest fascinerande när det gäller gallerna är att deras utseende ofta är helt olikt någon annan del av den växt som den bildas på. Dvs, gallerna är ofta stekelspecifika snarare än växtspecifika, och man kan ofta känna igen vilken art av stekel som har gjort gallen genom att titta på gallens utseende.

Gallsteklar tros ha sitt ursprung bland steklar som levde på vedlevande insektslarver, och de närmaste släktingarna till Cynipidae är alla sk parasitoider ("parasiter" som dödar sin värd). Inom familjen Cynipidae finns förutom gallbildande arter även steklar vars förfäder tros ha förlorat förmågan att bilda galler. Stekeln kan dock påverka den fortsatta utvecklingen av en redan påbörjad gall. Dessa steklar, även kallade "inhysingar", letar upp och lägger sina ägg i redan påbörjade eller färdiga galler. Till skillnad från parasitoider så äter inte inhysingarnas larver gallstekellarven utan livnär sig på växtmassan i gallen. I tillägg till inhysingarna så visar även andra insekter intresse för galler och dess innehåll. Äldre, tomma galler kan fungera som skydd för andra insekters ägg och larver, och gallbildar- och inhysingslarverna blir ofta parasiterade, inte sällan av andra stekelarter. Några av dessa parasitoida steklar (familjen Figitidae) tros vara gallsteklarnas närmaste släktingar. Utvecklingen av levnadssätt bland gallsteklarna och deras närmaste släktingar har således gått från att vara parasitoider, till att bli gallbildare, och till att förlora förmågan att bilda galler (inhysingar).

En tidigare hypotes om inhysingarnas ursprung var att de från början var gallbildare, men att de utvecklades till att bli inhysingar i galler som bildats av andra arter på samma typ av växt som de ursprungligen gjorde galler på.

Enligt den hypotesen skulle gallsteklar och inhysingar som lever på rosenbuskar vara varandras närmaste släktingar, istället för att tex rosgallsteklarna var närmare släkt med tex ekgallsteklarna. Senare studier visade att en mer trolig hypotes var att inhysingarna faktiskt hade ett unikt ursprung, och alla inhysingar var närmare släkt med varandra än de var med någon gallbildare. Dvs, inhysingarna ansågs vara en sk monofyletisk, eller naturlig grupp.

Tidigare studier av släktskap mellan gallsteklar har oftast baserats på morfologiska undersökningar. Denna avhandling innefattar den första större analysen av Cynipidae som baseras på molekylära data. Dessutom analyseras molekylära data tillsammans med tidigare publicerade morfologiska data i en sammanslagen släktskapsanalys. För att möjliggöra detta utvecklade vi statistiska metoder som baserar sig på bayesiansk statistik, s.k. betingade sannolikheter. De metoder vi använt uppskattar sannolikheten av en släktskaphypotes (fylogenetiskt träd, fylogeni), givet de DNA sekvenser och morfologiska karaktärer vi observerat. Manuskript I är en pilotstudie med ett fåtal gallstekelararter där vi undersöker hur olika DNA-sekvenser (gener) kan användas för att uppskatta gallsteklarnas fylogeni i en sk fylogenetisk analys. I manuskript II ökar vi antalet arter och använder några av de gener som vi funnit potentiellt användbara, och kombinerar dessa med data från detaljerade studier av gallsteklarnas morfologi. Vi visar även på hur nya statistiska metoder och modeller kan användas för att analysera DNA-sekvenser tillsammans med morfologiska karaktärer i fylogenetiska analyser. Manuskript III är en utvidgning av de teoretiska aspekterna av manuskript II, främst vad gäller hur man använder bayesianska metoder för att välja den analysmodell som har högst betingad sannolikhet. Det påpekas särskilt att det kan finnas osäkerhet i valet av en enskild modell, och att den osäkerheten i värsta fall kan påverka analyserna negativt. Det beskrivs hur man istället kan använda flera modeller i en fylogenetisk analys, och vikta var och en av modellerna i förhållande till dess betingade sannolikhet. Slutligen, i manuskript IV använder vi de data och metoder vi tidigare utvecklat för att göra en större analys baserad på 70 gallstekelararter tillsammans med ett 20 tal av de parasitiska släktingarna till Cynipidae.

De resultat från manuskript IV som enbart baserat sig på gendata (figur 3) var till stora delar likt de tidigare resultaten från morfologiska analyser (tidigare hypoteser sammanfattade i figur 1). Dock skiljde sig i ett antal hänseenden. En är grupperingen av en grupp gallsteklar som gör galler på vedartade växter. Dessa grupper, tillhörande triberna Pediaspidini, Diplolepidini, och Eschatocerini, placerade sig tillsammans i en förgrening nära basen på trädet, utanför resten av gallsteklarna (jämför figur 1 och figur 3). Den andra skillnaden jämfört med de morfologiska resultaten är att inhysingarna (tribus Synergini) visade sig inte vara varandras närmaste släktingar; de bildade således ingen monofyletisk grupp (se figur 3). Dessa

skillnader i resultat gör att slutsatserna vad det gäller gallsteklarnas evolution blir annorlunda. Med hjälp av bayesianska metoder härleddes de ursprungliga karaktärstillstånden för anfadern till gallsteklarna i familjen Cynipidae (resultaten redovisas i tabell 2). Det visade sig att de DNA-baserade resultaten stämde bra överens med tidigare resultat när det gällde slutsatsen att gallsteklarna troligtvis härstammar från palearktisk (främst nuvarande Europa), och att den första typen av galler troligen var distinkta utväxter, hade en kammare, och gjordes på frukter, frön, eller blommor. Däremot antyder resultaten från DNA-analyserna att de ursprungliga gallerna gjordes på vedartade bokväxter (Fagaceae) istället för på örtartade vallmoväxter (Papaveraceae).

Olika faktorer som kan bidra till dessa skillnader mellan resultaten från DNA och morfologi diskuteras i avhandlingen och i manuskript IV. Till exempel kan det hända att olika evolutionära linjer utvecklas med olika hastighet, och att de är alltför olika för att de modeller vi använder för att utröna släktskap skall klara av att hitta rätt släktskapsträd. Andra orsaker är att organismer som inte är närbesläktade kan leva i samma sorts miljö, och därför tendera att få samma utseende. Våra metoder för fylogenetisk analys kan förledas av dessa, sk konvergenta, likheter. Vi provade att använda komplexa modeller som är utvecklade för att hantera olikheter i evolutionär hastighet, men olikheterna bestod. Detta skulle antyda att en trolig anledning till att DNA och morfologi visar olika resultat är att gallsteklar med olika ursprung tenderar att se lika ut om de har liknande levnadssätt. Vi påpekar dock att det finns ytterligare potential för att utveckla mer sofistikerade analysmetoder, som skulle klara av att med större säkerhet peka på orsaker till de olikheterna vi sett. Kommande släktskapsanalyser av Cynipidae skulle dessutom underlättas om ytterligare data samlades in från gallsteklarna och deras släktingar.

Acknowledgements

I started studying biology because I wanted to be a scientist working with venomous snakes. I ended up working with worms, and then later with gall wasps. I still feel a little sad about the snakes but looking back I must admit that this journey among the phyla was the right way to go. It took me right into the heart of a field of science that I now have learned to love, namely systematics. Along this way, I have met people that inspired and taught me the art of the discipline, and without whom, I wouldn't be where I am, or who I am today. I wish to send my sincere thanks to those people. In (almost) chronological order they were: Göran Nilson, for showing me that there was systematics in the first place, Niklas Wikström, Torsten Eriksson and Hans-Erik Wanntorp, for showing me that by practicing phylogenetic systematics, I would join a "scientific revolution" (and that parts of this revolution could take place during discussions), Christer Erséus for his trust in me when he placed his molecular data in my hands, Mari Källersjö for her overwhelmingly enthusiasm and encouragements, Steve Farris and Pablo Goloboff for letting me take part of their brilliant thoughts (Pablo is also a strong climber, much stronger than me), Arnold Kluge, for reminding me of the philosophical part of systematics and science in general (although Arnold would probably say that I've been blinded by the light), John Huelsenbeck for providing much of this blinding light, Bengt Oxelman for most invaluable late-night discussions, wholehearted and sincere support (tack Bengt!), and for showing temper, and last but not least — actually the most — I wish to thank my supervisor Fredrik Ronquist. Without doubt, the most complete systematist I ever met, and I have felt the honor of being his student everyday since the first day he took me under his wings. Thank you Fredrik!

I also wish to thank the people that made the journey so pleasant: the lovely Afsaneh Ahmadzadeh, who carried me through the worst moments in the lab and helped me with her magic and without whom this thesis wouldn't be finished! Mikael Thollesson for allowing me to (ab)use his computers, and for fruitful and inspiring discussions (and for Perl!). All the students and staff I had the great pleasure to meet during my time at the EBC. Too many to mention but a few came closer than others; Per "Discokungen" Alström, Kristina "Kalenderflickan" Articus, Björn-Axel "Byrådirektörn" Beier, Per "Bootstrap" Erixon, Per "Huxley" Kornhall, Johan "Quake" Liljeblad, Anders "Masken" Lindström, Isabel "El Enciclopedia" Sanmartin, Magnus

"PUSS!" Popp, and Annika "Idolen" Vinnersten. Ni och ni andra, tack för all tid tillsammans!

An honour to Björn-Axel Beier, Jacob Höglund, Jobs-Karl Larsson, and Martin Irestedt, for appreciating what science is all about.

Mattias "Pennan" Forshage, and Bengt Oxelman gave valuable and most encouraging comments on an earlier version of the thesis summary, and Magnus Popp helped with preparing Figure 3.

Financial support was received from the Helge A:xon Johnson Foundation, and a from the Swedish Research Council (grant to Fredrik Ronquist).

The last years during the work on this thesis has for various reasons been the worst, but also the best years of my life. I want to return the love that got me going before, during, and after, to my family: Ingrid, Pål, Hanna, Pontus, Kim, Elvira, and Ebony. Finally, to one very special person:

"Kom med och gör en stund som skall vara kvar i hundra tusen år. Ibland vill man inte gärna gå in fast regnet öser ner. Ibland vill man hellre sitta kvar och skratta högt där döden ser Kom med och sätt ditt liv på spel för min skull."
(J. Hellman)

Tack Hege, för att du vågade!

References

- Akaike, H. 1973. Information theory as an extension of the maximum likelihood principle. Pages 267–281 *in* Second international symposium on information theory (B. N. Petrov and F. Csaki, eds.). Akademiai Kiado, Budapest.
- Angiosperm Phylogeny Group, 2003. An update of the angiosperm phylogeny group classification for the orders and families of flowering plants: APG II. *Bot. J. Lin. Soc.* 141:399–436.
- Archibald, J. K., M. E. Mort, and D. J. Crawford. 2003. Bayesian inference of phylogeny: a non-technical primer. *Taxon* 52:187–191.
- Askew, R. R. 1984. The biology of gallwasps. Pages 223–271 *in* Biology of gall insects (Anantakrishnan, T. N., ed.). Edward Arnold, London.
- Bollback, J. P. 2002. Bayesian model adequacy and choice in phylogenetics. *Mol. Biol. Evol.* 19:1171–1180.
- Bruno, W. J., and A. L. Halpern. 1999. Topological bias and inconsistency of maximum likelihood using wrong models. *Mol. Biol. Evol.* 16:564–566.
- Buckley, T. R., P. Arensburger, C. Simon, and G. K. Chambers. 2002. Combined data, Bayesian phylogenetics, and the origin of the New Zealand Cicada genera. *Syst. Biol.* 51:4–18.
- Burnham, K. P., and D. R. Anderson. 2002. Model selection and multimodel inference, a practical information-theoretic approach. Second edition. Springer, New York.
- Caterino, M. S., Cho, S., and F. A. H. Sperling. 2000. The current state of insect molecular systematics: a thriving tower of Babel. *Annu. Rev. Entomol.* 45:1–54.
- Drummond, A., and A. Rambaut. 2003. BEAST: Bayesian evolutionary analysis sampling trees. <http://evolve.zoo.ox.ac.uk/Beast/>.
- Edwards, A. W. F. 1992. Likelihood. Expanded edition. John Hopkins, London.
- Edwards, A. W. F., and L. L. Cavalli-Sforza. 1964. Reconstruction of evolutionary trees. Pages 67–76 *in* Phenetic and phylogenetic classification (Heywood, V. H., and J. McNeill, eds.). Systematics Association Publ. No. 6, London.
- Erixon, P. B., Svennblad, T. Britton, and B. Oxelman. 2003. Reliability of Bayesian posterior probabilities and bootstrap frequencies in phylogenetics. *Syst. Biol.* 52:665–673.
- Farris, J. S. 1999. Likelihood and inconsistency. *Cladistics* 15:199–204.
- Felsenstein, J. 1978. Cases in which parsimony or compatibility methods will be positively misleading. *Syst. Zool.* 27:401–410.
- Felsenstein, J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach. *J. Mol. Evol.* 17:368–376.
- Felsenstein, J. 1985. Confidence limits on phylogenies: an approach using the bootstrap. *Evolution* 39:783–791.
- Felsenstein, J. 2003. The troubled growth of statistical phylogenetics. *Syst. Biol.* 50:465–467.
- Felsenstein, J. 2003. Inferring phylogenies. Sinauer Associates, Sunderland, Massachusetts.

- Gauld, I., and B. Bolton. 1988. *The Hymenoptera*. Oxford University Press, Oxford.
- Gilks, W. R., S. Richardson, and D. J. Spiegelhalter. 1996. *Markov chain Monte Carlo in practice*. Chapman and Hall, London.
- Grant, T., and A. G. Kluge. 2003. Data exploration in phylogenetic inference: scientific, heuristic, or neither. *Cladistics*, 19:379–418.
- Hastings, W. 1970. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* 57:97–109.
- Holder, M. T. 2001. Using a complex model of sequence evolution to evaluate and improve phylogenetic methods. Ph. D. thesis, Univ. Texas at Austin, Austin.
- Holder, M., and P. O. Lewis. 2003. Phylogeny estimation: traditional and Bayesian approaches. *Nature Genetics* 4:275–284.
- Huelsenbeck, J. P. 1995. Performance of phylogenetic methods in simulation. *Syst. Biol.* 44:17–48.
- Huelsenbeck, J. P. 2002. Testing a covariotide model of DNA substitution. *Mol. Biol. Evol.* 19:698–707.
- Huelsenbeck, J. P., and K. A. Crandall. 1997. Phylogeny estimation and hypothesis testing using maximum likelihood. *Annu. Rev. Ecol. Syst.* 28:437–466.
- Huelsenbeck, J. P., and F. Ronquist. 2001. MrBayes: Bayesian inference of phylogeny. *Biometrics* 17:754–755.
- Huelsenbeck, J. P., and J. P. Bollback. 2001. Empirical and hierarchical Bayesian estimation of ancestral states. *Syst. Biol.* 50:351–366.
- Huelsenbeck, J. P., and N. S. Imenov. 2002. Geographic origin of human mitochondrial DNA: accommodating phylogenetic uncertainty and model comparison. *Syst. Biol.* 51:155–165.
- Huelsenbeck, J. P., B. Larget, R. E. Miller, and F. Ronquist. 2002. Potential applications and pitfalls of Bayesian inference of phylogeny. *Syst. Biol.* 51:673–688.
- Huelsenbeck, J. P., B. Rannala, and J. P. Masly. 2000. Accommodating phylogenetic uncertainty in evolutionary studies. *Science* 288:2349–2350.
- Huelsenbeck, J. P., F. Ronquist, R. Nielsen, and J. P. Bollback. 2001. Bayesian inference of phylogeny and its impact on evolutionary biology. *Science* 294:2310–2314.
- Kinsey, A. C. 1920. Phylogeny of cynipid genera and biological characteristics. *Bul. Am. Mus. Nat. Hist.* 42:357–402.
- Kishino, H., and H. Hasegawa. 1989. Evaluation of the maximum likelihood estimate of the evolutionary tree topologies from DNA sequence data, and the branching order in Hominoidea. *J. Mol. Evol.* 29:170–179.
- Larget, B., and D. L. Simon. 1999. Markov chain Monte Carlo algorithms for the Bayesian analysis of phylogenetic trees. *Mol. Biol. Evol.* 16:750–759.
- Lewis, P. O. 2001. Phylogenetic systematics turns over a new leaf. *Trends Ecol. Evol.* 16:30–37.
- Li, S. 1996. Phylogenetic tree construction using Markov chain Monte Carlo. Ph.D. thesis, Ohio State Univ., Columbus.
- Liljeblad, J. 2002. Phylogeny and evolution of gall wasps (Hymenoptera: Cynipidae). Ph.D. thesis, Stockholm Univ., Stockholm.
- Liljeblad, J., and F. Ronquist. 1998. A phylogenetic analysis of higher-level gall wasp relationships. *Syst. Entom.* 23:229–252.
- Madigan, D., and A. E. Raftery. 1994. Model selection and accounting for model selection uncertainty in graphical models using Occam's window. *J. Am. Stat. Assoc.* 89:1535–1546.
- Malyshev, S. I. 1968. *Genesis of the Hymenoptera and the phases of their evolution*. Methuen, London.

- Mau, B. 1996. Bayesian phylogenetic inference via Markov chain Monte Carlo. Ph.D. thesis, Univ. Wisconsin, Madison.
- Metropolis, N., A. Rosenbluth, M. Rosenbluth, A. Teller, and E. Teller. 1953. Equations of state calculations by fast computing machines. *J. Chem. Phys.* 21:1087–1092.
- Newton, M. A., and A. E. Raftery. 1994. Approximate Bayesian inference by the weighted likelihood bootstrap (with discussion). *J. Roy. Stat. Soc. B Met.* 56:3–48.
- Nieves-Aldrey, J. L. 2001. Hymenoptera, Cynipidae. *Fauna Ibérica*, Vol. 16. Mus. Nacion. Cienc. Natur., CSIC. Madrid.
- Posada, D., and K. A. Crandall. 2001. Selecting the best-fit model of nucleotide substitution. *Syst. Biol.* 50:580–601.
- Pupko, T., D. Huchon, Y. Cao, N. Okada, and M. Hasegawa. 2002. Combining multiple data sets in a likelihood analysis: which models are the best? *Mol. Biol. Evol.* 19:2294–2307.
- Rasnitsyn, A. P. 1988. An outline of evolution of the hymenopterous insects. *Orient. Insects* 22:115–145.
- Rogers, J. S. 2001. Maximum likelihood estimation of phylogenetic trees is consistent when substitution rates vary according to the invariable sites plus gamma distribution. *Syst. Biol.* 50:713–722.
- Ronquist, F. 1994. Evolution of parasitism among closely related species: phylogenetic relationships and the origin of inquilism in gall wasps (Hymenoptera: Cynipidae). *Evolution* 48:241–266.
- Ronquist, F. 1995. Phylogeny and early evolution of the Cynipoidea (Hymenoptera). *Syst. Entom.* 20:309–335.
- Ronquist, F. 1999. Phylogeny, classification and evolution of the Cynipoidea. *Zool. Scripta* 28:139–164.
- Ronquist, F., and J. P. Huelsenbeck. 2003. MRBAYES 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics* 19:1572–1574.
- Ronquist, F., and J. P. Huelsenbeck. in press. A Bayesian approach to supertrees. *In Phylogenetic supertrees: Combining information to reveal the tree of life* (O. R. P. Bininda-Emonds, ed.). Kluwer Academic, Dordrecht.
- Ronquist, F., and J. Liljeblad. 2001. Evolution of the gall wasp-host plant association. *Evolution* 51:2503–2522.
- Ronquist, F., A. P. Rasnitsyn, A. Roy, K. Eriksson, and M. Lindgren. 1999. Phylogeny of the Hymenoptera: A cladistic reanalysis of Rasnitsyn's (1988) data. *Zool. Scripta* 28:13–50.
- Sanderson, M. J., and J. Kim. 2000. Parametric phylogenetics? *Syst. Biol.* 49:817–829.
- Schultz, T. R., and G. A. Churchill. 1999. The Role of subjectivity in reconstructing ancestral character states: a Bayesian approach to unknown rates, states, and transformation asymmetries. *Syst. Biol.* 48:651–664.
- Schwartz, G. 1978. Estimating the dimensions of a model. *Ann. Stat.* 6:461–464.
- Steel, M., and D. Penny. 2000. Parsimony, likelihood, and the role of models in molecular phylogenetics. *Mol. Biol. Evol.* 17:839–850.
- Sullivan, J., and D. L. Swofford. 2001. Should we use model-based methods for phylogenetic inference when we know that assumptions about among-site rate variation and nucleotide substitution pattern are violated? *Syst. Biol.* 50:723–729.
- Suzuki, Y., G. V. Glazko, and M. Nei. 2002. Overcredibility of molecular phylogenetics obtained by Bayesian phylogenetics. *Proc. Natl. Acad. Sci. USA* 99:16138–16143.

- Swofford, D. L. 2002. PAUP*. Phylogenetic analysis using parsimony (* and other methods). Version 4b10. Sinauer Associates, Sunderland, Massachusetts.
- Swofford, D. L., P. J. Waddell, J. P. Huelsenbeck, P. G. Foster, P. O. Lewis, and J. S. Rogers. 2001. Bias in phylogenetic estimation and its relevance to the choice between parsimony and likelihood methods. *Syst. Biol.* 50:525–539.
- Tierney, L. 1994. Markov chains for exploring posterior distributions (with discussion). *Ann. Stat.* 22:1701–1762.
- Tuffley, C., and M. Steel. 1998. Modeling the covarion hypothesis of nucleotide substitution. *Math. Biosci.* 147:63–91.
- Vårdal, H. 2004. From parasitoids to gall inducers and inquilines: morphological evolution in gall wasps. Ph. D. Thesis, Univ. Uppsala, Uppsala.
- Wasserman, L. 2000. Bayesian model selection and model averaging. *J. Math. Psych.* 44:92–107.
- Whelan, S., P. Liò, and N. Goldman. 2001. Molecular phylogenetics: state-of-the-art methods for looking into the past. *Trends Genet.* 17:262–272.
- Wilcox, T. P., D. J. Zwickl, T. A. Heath, and D. M. Hillis. 2002. Phylogenetic relationships of the dwarf boas and a comparison of Bayesian and bootstrap measures of phylogenetic support. *Mol. Phyl. Evol.* 25:361–371.
- Yang, Z., and B. Rannala. 1997. Bayesian phylogenetic inference using DNA sequences: a Markov chain Monte Carlo method. *Mol. Biol. Evol.* 14:717–724.
- Zucchini, W. 2000. An introduction to model selection. *J. Math. Psych.* 44:41–61.

Acta Universitatis Upsaliensis

*Comprehensive Summaries of Uppsala Dissertations
from the Faculty of Science and Technology*

Editor: The Dean of the Faculty of Science and Technology

A doctoral dissertation from the Faculty of Science and Technology, Uppsala University, is usually a summary of a number of papers. A few copies of the complete dissertation are kept at major Swedish research libraries, while the summary alone is distributed internationally through the series *Comprehensive Summaries of Uppsala Dissertations from the Faculty of Science and Technology*. (Prior to October, 1993, the series was published under the title “Comprehensive Summaries of Uppsala Dissertations from the Faculty of Science”.)

Distribution:

Uppsala University Library
Box 510, SE-751 20 Uppsala, Sweden
www.uu.se, acta@ub.uu.se

ISSN 1104-232X
ISBN 91-554-5872-6