RH: Phylogenetic analysis of host repertoire evolution

# Bayesian inference of ancestral host-parasite interactions under a phylogenetic model of host repertoire evolution

Mariana P Braga[1], Michael Landis[2], Sören Nylin[1], Niklas Janz[1] and Fredrik Ronquist[3]

[1]*Department of Zoology, Stockholm University, Stockholm, SE-10691, Sweden;*

[2]*Department of Ecology and Evolution, Yale University, New Haven, CT, 06511, USA;*

[3]*Department of Bioinformatics and Genetics, Swedish Museum of Natural History, Stockholm, Sweden*

**Corresponding author:** Mariana P Braga, Department of Zoology, Stockholm University, Stockholm, SE-10691, Sweden; E-mail: mariana.braga@zoologi.su.se

1  *Abstract.*— Intimate ecological interactions, such as those between parasites and their hosts, may persist

2  over long time spans, coupling the evolutionary histories of the lineages involved. Most methods that

3  reconstruct the coevolutionary history of such associations make the simplifying assumption that parasites

4  have a single host. Many methods also focus on congruence between host and parasite phylogenies, using

5  cospeciation as the null model. However, there is an increasing body of evidence suggesting that the host

6  ranges of parasites are more complex: that host ranges often include more than one host and evolve via

7  gains and losses of hosts rather than through cospeciation alone. Here, we develop a Bayesian approach for

8  inferring coevolutionary history based on a model accommodating these complexities. Specifically, a

9  parasite is assumed to have a host repertoire, which includes both potential hosts and one or more actual

10  hosts. Over time, potential hosts can be added or lost, and potential hosts can develop into actual hosts or

11  vice versa. Thus, host colonization is modeled as a two-step process, which may potentially be influenced

12  by host relatedness or host traits. We first explore the statistical behavior of our model by simulating

13  evolution of host-parasite interactions under a range of parameters. We then use our approach,

14  implemented in the program RevBayes, to infer the coevolutionary history between 34 Nymphalini

15  butterfly species and 25 angiosperm families.

16  (Keywords: ancestral hosts, coevolution, herbivorous insects, probabilistic modeling.)

17    Extant ecological interactions, such as those between parasites and hosts, are often the

18  result of a long history of coevolution between the involved lineages (Elton 1946; Klassen 1992).

19  Specialization is predominant among parasites (including parasitic herbivorous insects; Forister

20  et al. 2015), but host associations are not static: they continuously evolve over time via gains and

21  losses of hosts (Janz and Nylin 2008; Nylin et al. 2018). The colonization of new hosts and loss of

22  old hosts not only shape the evolutionary trajectories of the interacting lineages, but can also

23  have large effects at ecological timescales (Nosil 2002; Calatayud et al. 2016). These effects are

24  evident, for example, with emerging infectious diseases and zoonotic diseases (Acha and Szyfres

25  2003), which involve colonization of new hosts within and among groups of domesticated species

26  (Subbarao et al. 1998), wildlife (Fisher et al. 2009), and humans (Hahn et al. 2000). Unraveling

27  the processes underlying changes in species associations is thus key to understanding evolutionary

28  and ecological phenomena at various timescales, such as the emergence of infectious diseases,

29  community assembly, and parasite diversification (Hoberg and Brooks 2015).

30    Many methods developed to study historical associations focus on congruence between

31  host and parasite phylogenies (Brooks 1979; Huelsenbeck et al. 1997; de Vienne et al. 2013). Such

32  methods largely fall into two main classes of cophylogenetic approaches: (1) topology- and

33  distance-based methods, which estimate the congruence between two phylogenies (Legendre et al.

34  2002), and (2) event-based methods, which map the parasite phylogeny onto the host phylogeny

35  using evolutionary events (Ronquist 2003). Typically, cospeciation is the null hypothesis in these

36  methods, where host shifts are invoked only to explain deviations from cospeciation (de Vienne

37  et al. 2013). Moreover, most of these methods do not allow ancestral parasites to be associated

38  with more than one host lineage, thus failing to account for a potentially important driver of

39  parasite diversification (Janz and Nylin 2008).

40    An alternative approach to studying coevolving host-parasite associations is to perform

41  ancestral state reconstructions of individual host taxa onto the parasite phylogeny and combine

42  the ancestral host states *a posteriori* into inferred host ranges (e.g. Nylin et al. 2014). Even

43  though this approach allows ancestral parasites to have multiple hosts, it assumes that the

44  associations between the parasite and each host evolve independently. This has a number of

45  serious drawbacks. For instance, ancestral parasites may be inferred to have an unrealistically

46  high number of hosts, or no host at all. Furthermore, the more narrowly circumscribed the host

47  taxa are, the more likely it is that ancestral parasite lineages are reconstructed as having no

48  hosts. In addition, the independence assumption causes the phylogenetic relationships among

49  hosts to be ignored, meaning that the model assigns equal rates to all colonizations of new hosts

50  regardless of how closely related the new host is to the current hosts being used by the parasite.

51      A desirable model of host usage should therefore allow parasites to have multiple hosts,

52  while also allowing for among-host (or context-dependent) effects to influence ancestral host use

53  estimates and gain and loss rates in whatever manner explains the biological data best. One

54  possible solution is to restate the problem of host-parasite co-evolution in terms of historical

55  biogeography. For instance, the Dispersal-Extirpation-Cladogenesis (DEC) model of Ree et al.

56  (2005) allows species ranges to stochastically evolve as a set of discrete areas over time through

57  area gain events (dispersal), area loss events (extirpation), and cladogenetic events (range

58  inheritance patterns that reflect speciational models). Although these methods are designed for

59  biogeographic inference, a similar approach is clearly suitable for more realistic modeling of

60  host-parasite coevolution dynamics, where colonization and loss of hosts (instead of discrete

61  areas) is modeled as a continuous-time Markov process (e.g. Hardy 2017). In biogeography, the

62  colonization of a new area or the disappearance from a previously occupied area is modeled as a

63  binary trait: the species is either present or absent in the area. While this binary view might be

64  simple but useful in biogeography, it may be too simplistic for use in the coevolution between

65  hosts and parasites. For instance, it is known that butterflies can utilize a range of plants that

66  they do not regularly feed on in the wild, and it has been suggested that these potential hosts

67  have played an important role in the evolution of host use in butterflies, by increasing the

68  variability in host use through time and across clades (Janz et al. 2016; Braga et al. 2018). This

69  hypothesis can only be directly tested, however, if we explicitly model the evolution of host use as

70  a two-step process, which cannot be done with the binary methods that are used today to study

71  host-parasite coevolution or biogeography.

72      Here, we propose a model where a parasite is assumed to have a *host repertoire*, defined as

73  the set of all potential and actual hosts for that parasite. In this model, the colonization of a new

74  host involves two steps: first, the parasite gains the ability to use the new host (it becomes a

75  potential host), and then starts actually using it in nature (it becomes an actual host). These two

76  steps can be interpreted as the inclusion of the new host into the fundamental and then into the

77  realized host repertoire of the parasite - analogous to fundamental and realized niche (Nylin et al.

78  2018; Larose et al. 2019). Similarly, the complete loss of a host from a parasite's realized

79  repertoire involves two steps. First, it changes from an actual to a potential host, and then it is

80  lost completely from the host repertoire. For example, if the geographic range of a host

81  contracted to become allopatric with respect to a parasite's geographic range, the host would

82  remain as part of the fundamental repertoire until the parasite completely lost the ability to use

83  the host, in which case the host would be lost from the repertoire. Even when in sympatry, the

84  evolution of a new defense mechanism by the host may prevent the parasite from using that host.

85  However, since host use is a complex and multidimensional trait, it is unlikely that a parasite

86  loses all the machinery necessary to use a host in one single event, and it may well retain some

87  ability to survive on the host. Thus, three host-parasite association states are necessary for such a

88  two-step model: the host is used (actual host), the parasite has some ability to use the host but

89  does not use it in nature (potential host), and the parasite cannot use the host (non-host).

90      In this paper, we develop a Bayesian approach to coevolutionary inference based on such a

91  model of host repertoire evolution, inspired by the previous work on similar biogeographic

92  inference problems by Landis et al. (2013). The basic binary biogeographic model, when applied

93  to coevolution, accommodates both multiple ancestral hosts and changes in host configurations

94  over time that correspond to evolutionary changes in host lineages or host traits. We extend this

95  model to also include a two-step host colonization process, such that the fundamental host

96  repertoire can persist over time and affect the evolution of the realized repertoire. We have

97  implemented the model in RevBayes (Höhna et al. 2016), allowing us to perform simulation as

well as Bayesian Markov chain Monte Carlo (MCMC) inference under the model. This Bayesian framework allows one to estimate the joint distribution of host gain and loss rates, the effect (if any) of phylogenetic distances among hosts upon host gain rates, and the historical sequences of evolving host repertoires among the parasites. Using simulations, we explore the statistical behavior of our approach, and demonstrate its empirical application with an analysis of the coevolution between Nymphalini butterflies and their angiosperm hosts.

# Methods

## *Model description*

We are interested in modeling the evolution of ecological interactions between $M$ extant parasite taxa and $N$ host taxa, where each parasite uses one or more hosts. Rooted and time-calibrated phylogenetic trees describe the evolutionary relationships among the $M$ parasite taxa and among the $N$ host taxa. In this study, the trees are considered to be known without error. In principle, it would be straightforward for the model to accommodate phylogenetic uncertainty in the host or parasite trees but MCMC inference may prove challenging under such conditions.

Each parasite taxon has a host repertoire, which is represented by a vector of length $N$ that contains the information about which hosts the given parasite uses. The interaction between the $m$-th parasite and the $n$-th host is denoted $x_{m,n}$. At any given time, each host taxon can assume one of three states with respect to a parasite lineage: $x_{m,n}$ is equal to 0 (non-host), 1 (potential host), or 2 (actual host). Criteria for how to code non-host, potential host, and actual host states will depend on the host-parasite system under study; below, we provide criteria for our Nymphalini dataset that may act as guidelines. We allow all host repertoires in which the parasite has at least one actual host. Thus, the state space, $S$, includes $3^N - 2^N$ host repertoires for $N$ hosts.

122        Here we define the transition from state 0 to state 1 as the gain of the ability to use the

123  host, and the transition from state 1 to state 2 as the time when the parasite actually starts to

124  use the host in nature. If we assume that gains and losses of hosts occur according to a

125  continuous-time Markov chain, the probability of a given history of association between a parasite

126  clade and their hosts can be easily calculated (Ree and Smith 2008). This calculation is based on a

127  matrix, $\mathbf{Q}$, containing the instantaneous rates of change between all pairs of host repertoires, and

128  thus describing the Markov chain. Based on the $\mathbf{Q}$ matrix, it is possible to calculate the transition

129  probability of the observed host repertoires at the tips of the parasite tree by marginalizing over

130  the infinite number of histories that could produce the observed host repertoires. Unfortunately,

131  computing these transition probabilities becomes intractable as the number of host repertoire

132  configurations, $S$, grows large. Modeling host repertoire evolution for host repertoire size $N = 7$

133  requires an $S \times S$ rate matrix defined for $S = 3^7 - 2^7 = 2059$, causing $\mathbf{Q}$ to be too large for

134  efficient inference. In order to handle large host repertoires, we numerically integrate over possible

135  histories using data augmentation and MCMC rather than analytically computing the

136  probabilities using matrix exponentiation. This data augmentation approach has been used to

137  model sequence evolution for protein-coding genes (Robinson 2003) and historical biogeography

138  (Landis et al. 2013; Quintero and Landis 2019), suggesting the framework may be useful to model

139  host-parasite interactions as well. In this study, we assume that both daughter lineages identically

140  inherits their host repertoires from their immediate ancestor at the time of cladogenesis.

141        We define a model where the gain of a host (both 0→1 and 1→2) depends on the

142  phylogenetic distance between the available hosts and those currently used by a lineage. Figure 1

143  schematically illustrates the evolutionary dynamics of the model using $M = 4$ parasite species and

144  $N = 5$ host species, while assuming that host gain rates are independent (Fig. 1a,c) or dependent

145  (Fig. 1b,d) of phylogenetic distances among hosts. To formalize these dynamics, let $q_{\mathbf{y},\mathbf{z}}^{(a)}$ be the

146  rate of change from host repertoire $\mathbf{y}$ to repertoire $\mathbf{z}$ by changing the state of host $a$. Also, let $\lambda_{ij}$

147  be the rate at which an individual host changes from state $i$ to state $j$, and $\eta(\mathbf{y}, a, \beta)$ be a

148  phylogenetic-distance rate modifier. The phylogenetic-distance rate modifier function, $\eta$, rescales

149 the base rate of host gain to allow new hosts that are closely related to the parasite's current

150 hosts to be colonized at higher rates than distantly related hosts. We define the instantaneous

151 rate of change as

$$
q_{\mathbf{y},\mathbf{z}}^{(a)} = \begin{cases} \lambda_{10}, & \text{if potential host loss } (y_a = 1 \text{ and } z_a = 0) \\ \lambda_{01}\eta_1(\mathbf{y}, a, \beta) & \text{if potential host gain } (y_a = 0 \text{ and } z_a = 1) \\ \lambda_{21}, & \text{if actual host loss } (y_a = 2 \text{ and } z_a = 1) \\ \lambda_{12}\eta_2(\mathbf{y}, a, \beta) & \text{if actual host gain } (y_a = 1 \text{ and } z_a = 2) \\ 0, & \text{if direct transition between states 0 and 2 } (|y_a - z_a| > 1) \\ 0, & \text{if } \mathbf{y} \text{ and } \mathbf{z} \text{ differ at more than one host} \\ 0 & \text{if } \mathbf{z} \text{ does not contain at least one actual host} \end{cases}
$$

152 and the phylogenetic-distance rate modifier function as

$$
\eta(\mathbf{y}, a, \beta) = e^{-\beta d/\overline{d}}, \tag{1}
$$

153 where $\beta$ controls the effect of $d$, the average pairwise phylogenetic distance between the new host,

154 $a$, and the hosts currently occupied in $\mathbf{y}$; and $\overline{d}$ is the average phylogenetic distance between all

155 pairs of hosts. Pairwise phylogenetic distance is defined as the sum of branch lengths separating

156 two leaf nodes. The difference between $\eta_1$ and $\eta_2$ is that in the first, pairwise distances are

157 calculated between the new host and all potential and actual hosts, while in the second only

158 actual hosts are included. This allows for a model formulation where the effect of host distances

159 on $\lambda_{01}$ and on $\lambda_{12}$ are independent, while still allowing a formulation where they are equal. If

160 $\beta = 0$, the gain rate of host $a$ is equal to the unmodified gain rate, $\lambda_{01}$ or $\lambda_{12}$. If $\beta > 0$, the gain

161 rate of phylogenetically close hosts is higher than distant hosts.

162       We fit this model using the Bayesian data augmentation strategy described in Landis

163  et al. (2013). The method estimates the joint posterior probability of model parameters,

164  $\theta = (\mu, \lambda, \beta)$, and data-augmented evolutionary histories, $X_{\mathrm{aug}}$, conditional on the observed host

165  repertoire data, $X_{\mathrm{obs}}$, and the parasite phylogeny, $\Psi_{\mathrm{p}}$, and the host phylogeny, $\Psi_{\mathrm{h}}$, using MCMC.

166  To sample values from the posterior, $P(X_{\mathrm{aug}}, \theta \mid X_{\mathrm{obs}}, \Psi_p, \Psi_h)$, new parameter values for $\mu$, $\lambda$, and

167  $\beta$ are proposed using standard Metropolis-Hastings proposals for updating simple parameters

168  (Hastings 1970). Analogously, our MCMC stochastically proposes and/or accepts new augmented

169  host repertoire histories using the Metropolis-Hastings algorithm. Augmented histories are

170  proposed using two types of MCMC moves: branch-specific moves and node-and-branch moves.

171  Branch-specific moves propose a new augmented history by sampling a branch from the

172  phylogeny uniformly at random, then proposing new histories for a subset of host-characters using

173  the rejection sampling method of Nielsen (2002) under the assumption that all host characters

174  evolved under mutual independence ($\beta = 0$); this assumption allows us to rapidly propose new

175  augmented histories. Although augmented histories are proposed assuming host characters evolve

176  independently, we compute the acceptance probability for the branch-specific move by considering

177  the full-featured model probability that allows for non-independent rates of character change

178  when calculating the Metropolis-Hastings ratio. Thus, the augmented histories are sampled in

179  proportion to their posterior probabilities under the full model. Node-and-branch moves involves

180  sampling new host repertoire states for a node sampled uniformly at random within the parasite

181  tree, along with the three branches incident to the node. Together, the branch-specific moves, the

182  node-and-branch moves, and the parameter moves allow MCMC to estimate the posterior

183  probability of combinations of host repertoire histories and evolutionary parameters. Further

184  details are provided in Landis et al. (2013).


185                                        *Model selection*


186  When $\beta = 0$, the phylogenetic-distance dependent model, $M_D$ becomes a mutual-independence

187  model, $M_0$, where the interaction between the parasite and each host evolves independently.

188  These models are therefore nested ($M_0 \subseteq M_D$) and we can compute Bayes factors for model $M_D$
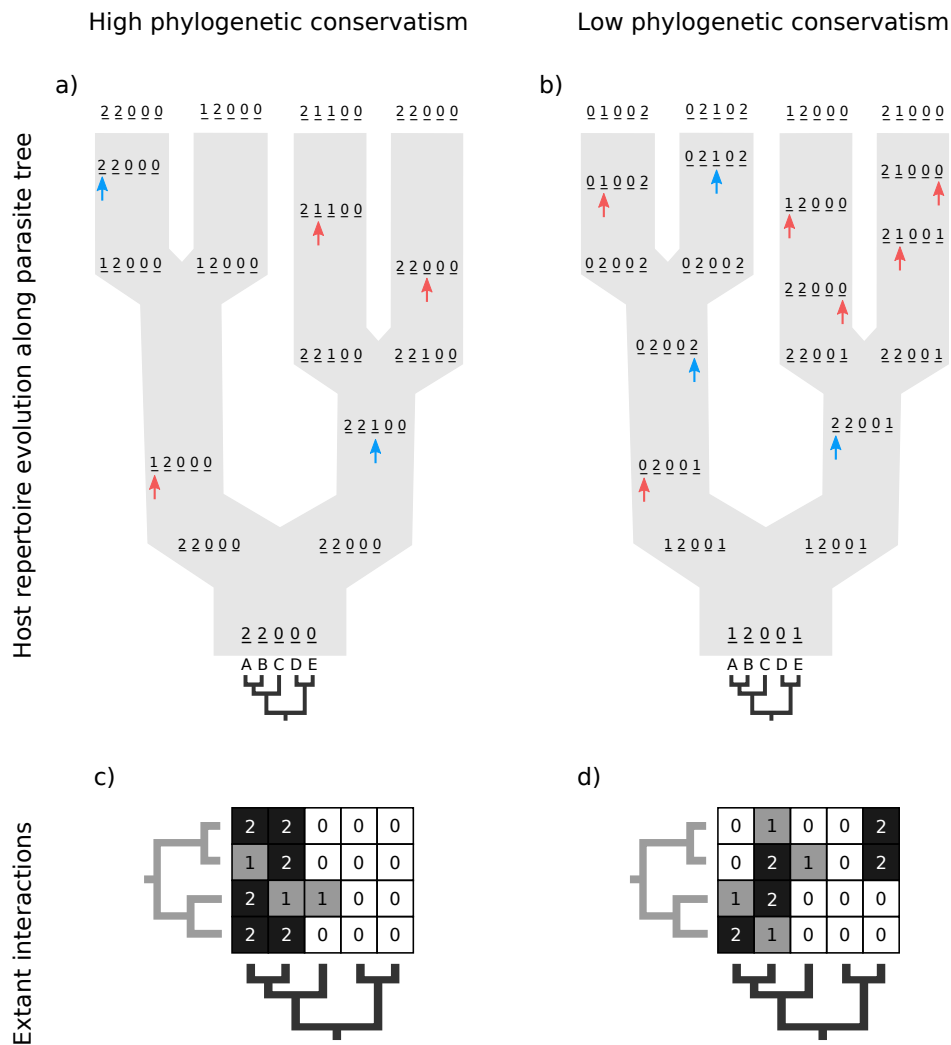
Figure 1: Host repertoire evolution along a hypothetical tree and resulting host-parasite interactions. Two examples of coevolutionary histories between four parasites and five hosts are shown to illustrate how the model works. Host repertoires evolve by gains ($0{\rightarrow}1$ and $1{\rightarrow}2$, blue arrows) and losses ($1{\rightarrow}0$ and $2{\rightarrow}1$, red arrows). Coevolutionary histories in **a** and **b** produce the interactions in **c** and **d** respectively. In **c** and **d**, each column represents one host and each row represents the host repertoire of one parasite. High phylogenetic conservatism is produced when the rate of repertoire evolution, $\mu$, is low and the effect of the phylogenetic distance between hosts, $\beta$, is high. Conversely, low phylogenetic conservatism is produced when $\mu$ is high and $\beta$ is low.

over model $M_0$ using the Savage-Dickey ratio (Verdinelli and Wasserman 1995; Suchard et al.

2001), defined as

$$B_{D,0} = \frac{P(\beta = 0 \mid M_D)}{P(\beta = 0 \mid \mathbf{x}_{\text{obs}}, M_D)} \tag{2}$$

where $P_D(\beta = 0 \mid M_D)$ is the prior probability and $P(\beta = 0 \mid \mathbf{x}_{\text{obs}}, M_D)$ is the posterior

probability, both defined in terms of the phylogenetic-distance dependent model, $M_D$, at the

restriction point $\beta = 0$ where $M_D$ and $M_0$ are equivalent. While we could directly compute the

prior probability of $\beta = 0$, we approximated the posterior at $\beta = 0$ using a kernel density

estimator with a gamma function, which only takes positive values, and a bandwidth of 0.02. To

interpret if and how Bayes factors favored the phylogenetic-distance dependent model, $M_D$, we

followed the guidelines of Jeffreys (1961): model $M_0$ is favored for Bayes factors with values less

than 1, insubstantial support is awarded to model $M_D$ for values between 1 and 3, substantial

support for values between 3 and 10, strong support for values between 10 and 30, very strong

support for values between 30 and 100, and decisive support for values greater than 100.

## *Data analysis*

*Simulation study.*— We simulated 50 datasets for each of nine combinations of values for the rate

of host-repertoire evolution, $\mu$ (0.01, 0.04, and 0.1), and values of $\beta$ (0, 1, and 4). These

parameter combinations produce datasets with varying degrees of phylogenetic conservatism for

both parasites and hosts (Fig. 2). Each dataset contained 34 insects and 25 hosts, and was

produced by simulating host repertoire evolution in the parasite tree used in the empirical study

(see below). Host gain and loss rates were chosen to resemble the rates inferred from the

empirical analysis. This simulation was designed to assess our statistical power to detect the

effect of phylogenetic distance among hosts upon host gain rates given the size of our empirical

dataset and the type of variation we expected it to contain.

We ran independent MCMC analyses for each set of 50 datasets, under the

phylogenetic-distance dependent model. We then quantified how well the posterior probabilities
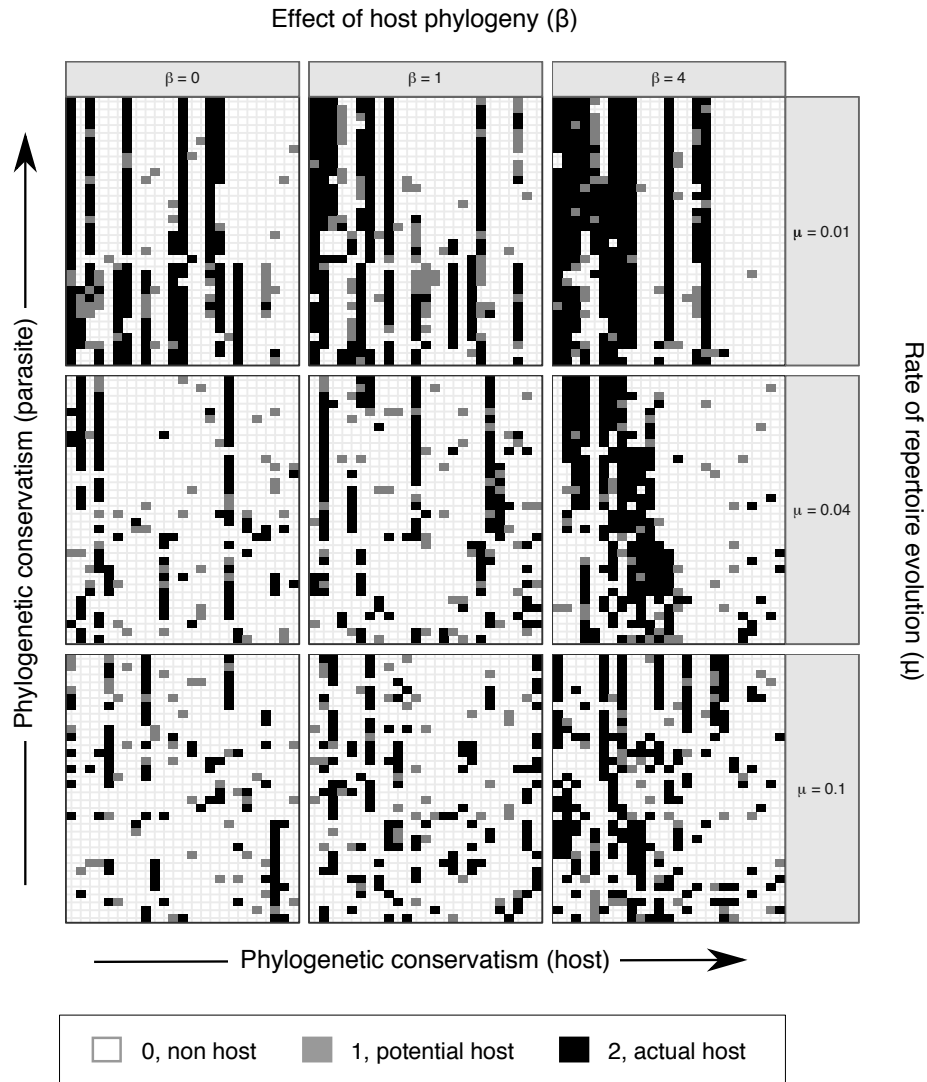
Figure 2: Simulated datasets for nine parameter combinations. Interactions between Nymphalini butterflies and their host plants for one of 50 simulations with each parameter combination. In each of the nine datasets, each column represents one host in the repertoire and each row shows the host repertoire of one butterfly species. When phylogenetic conservatism in host-parasite interactions is low for both hosts and parasites, the interactions are more randomly spread (matrix at bottom-left corner). As phylogenetic conservatism among parasites increases, host repertoires (rows) become more similar (upper matrices). When there is phylogenetic conservatism among hosts, host repertoires include more closely-related hosts (neighbouring columns; matrices to the right)

213 of coevolutionary histories correspond to the true history known from each simulation.

214 Specifically, we first computed the posterior probability of interaction between each host and each

215 internal node in the butterfly tree, for states 1 and 2 separately. Then, we calculated the sum of

216 squared differences between each posterior probability ($0 \leq P \leq 1$) and the corresponding truth

217 for that simulation (1, if the host was on the given state in the simulated dataset; 0, if not). This

218 error term increases as the inferred ancestral host repertoires become less accurate.

219 *Empirical study.*— In order to validate our method, we compiled data from the literature for

220 butterflies from the tribe Nymphalini (Nymphalidae) and their host plants (see Supplementary

221 Information for reference list). We chose this butterfly clade because we expect that a large

222 fraction of the real potential hosts are known, as there has been systematic experimental studies

223 of larval feeding ability. The dataset included 34 butterflies species and plants from 16

224 angiosperm families (Figs. S1 and S2). For each butterfly species, host plants commonly used in

225 nature were coded as 'actual hosts' and plants never used were coded as 'non-hosts'. Plants that

226 are not commonly used in nature, but for which there is strong evidence (field observation or

227 experiment) that the larvae can feed upon them, were coded as 'potential hosts'.

228 Because we lack the information on potential hosts for most host-parasite systems (i.e.

229 hosts are usually only classified as hosts or non-hosts), we tested whether our model is able to

230 recover the same parameter estimates and coevolutionary histories when all the potential hosts

231 are coded as non-hosts. For that, we ran the same analysis as for the full dataset, but first

232 removed all the 1s from the empirical dataset. Then we compared the posterior probabilities

233 inferred from the two datasets. To assess the similarities between the coevolutionary histories

234 inferred using the different datasets, we calculated summary statistics for the absolute difference

235 in probability of each interaction between hosts and internal nodes in the butterfly tree.

236 For both the simulation and empirical studies we used the phylogenetic relationships

237 between butterfly species in the Nymphalini tribe as proposed by Chazot (unpublished, Fig. S3)

238 and the phylogenetic relationships between angiosperm families proposed by Magallón et al.

239 (2015). Although our framework allows the inclusion of a large number of hosts in the same

240 analysis, computational time increases significantly with the size of the host repertoire. We

241 therefore chose to include 25 hosts, which allows the inclusion of all host lineages used by any of

242 the butterflies. To ensure the inclusion of all plant lineages that might have been used as hosts in

243 the past, we pruned the angiosperm phylogenetic tree so that all 16 families in the dataset were

244 included, and the remaining branches were collapsed to more ancestral nodes until only 25 tips

245 were left. We then pruned all the branches leading up to the tips to the time of origin of the

246 butterfly clade (approx. 22 Ma), and this pruned tree was then used to calculate phylogenetic

247 distances between hosts. To simplify the analysis, we hold the phylogenetic distances between

248 plant families constant, independent of geological time, even though the distances would be

249 expected to increase as evolution proceeds towards the recent.

250 We summarized inferred coevolutionary histories in two ways. First, we calculated the

251 posterior probability for fundamental and realized host repertoires at internal nodes of the

252 Nymphalini phylogeny based on the frequency with which states 1 and 2 were sampled for each

253 host during MCMC. Second, in order to reduce the dimensionality of the host repertoire and

254 facilitate visualization of ancestral state reconstructions, we assigned hosts to modules based on

255 extant butterfly-plant interactions (Fig. S2). Modules are groups of plants and butterflies that

256 interact more with each other than with other taxa, thus host plants are assigned to the same

257 module when they are used by the same butterflies. To identify the modules, we used a simulated

258 annealing algorithm that maximizes the index of modularity. Specifically, we used Newman and

259 Girvans metric (Newman and Girvan 2004) modified for bipartite networks (Barber 2007) as

260 implemented in the software MODULAR (Marquitti et al. 2014).

261 *Software configuration.*— All analyses were performed in RevBayes (Höhna et al. 2016). For the

262 simulated data, we ran two independent MCMC analyses for $10^5$ cycles, sampling parameters and

263 node histories every 50 cycles, and discarding the first $5 \times 10^4$ as burnin. For the empirical data,

264 we ran five independent MCMC analyses, each set to run for $10^6$ cycles, sampling every 50 cycles,

265 and discarding the first $10^5$ as burnin. To verify that MCMC analyses converged to the same

266 posterior distribution, we applied the Gelman diagnostic (Gelman and Rubin 1992) provided

267  through the R package *coda* (Plummer et al. 2006). For both simulated and empirical datasets,

268  we used the following priors: $\beta \sim \text{Exponential}(1)$, $\mu \sim \text{Exponential}(10)$, and

269  $\lambda \sim \text{Dirichlet}(1, 1, 1, 1)$. Analysis scripts and data files are available at

270  https://github.com/mpiresbr/host_repertoire. A RevBayes tutorial for the empirical

271  analysis will be soon available at https://revbayes.github.io/tutorials#host_rep.

## Results

273  *Simulation study.*—Posterior distributions of parameter values for the $9\times100$ MCMC analyses are

274  shown in Figure 3. Overall, the model was able to accurately recover the true simulation

275  parameters (true value within 95% highest posterior density, or HPD). However, accuracy

276  decreased with increasing rate of host repertoire evolution, possibly due to character saturation.

277  We performed model selection based on Bayes factors. Considering that the prior

278  distribution is $\beta \sim Exponential(1)$, a high marginal posterior probability for $\beta = 0$ under $M_D$ is

279  necessary to result in a Bayes factor $< 1$ and thus selection of $M_0$. For simulations with $\beta = 0$,

280  the correct model, $M_0$, was selected in more than 60% of the simulations, and most of the

281  remaining simulations gave insubstantial support to $M_D$ (Fig. 4). When $\beta = 1$, Bayes factors

282  correctly selected $M_D$ in the majority of cases, but strong support for $M_D$ was only achieved in

283  simulations with $\beta = 4$, particularly when the rate of evolution was highest ($\mu = 0.1$).

284  We then compared the true coevolutionary history of each simulation to the corresponding

285  posterior distribution of the sampled coevolutionary histories (Fig. 5). The estimation error, that

286  is, the sum of squared differences between estimated and true coevolutionary histories, was very

287  low when the rate of host-repertoire evolution was lowest ($\mu = 0.01$), but also when the

288  phylogenetic-distance power was highest ($\beta = 4$). This means that accuracy in the estimation of

289  coevolutionary history increases with phylogenetic conservatism on both the butterfly and the

290  plant trees. Overall, error was higher on the estimation of actual hosts (state 2) than potential
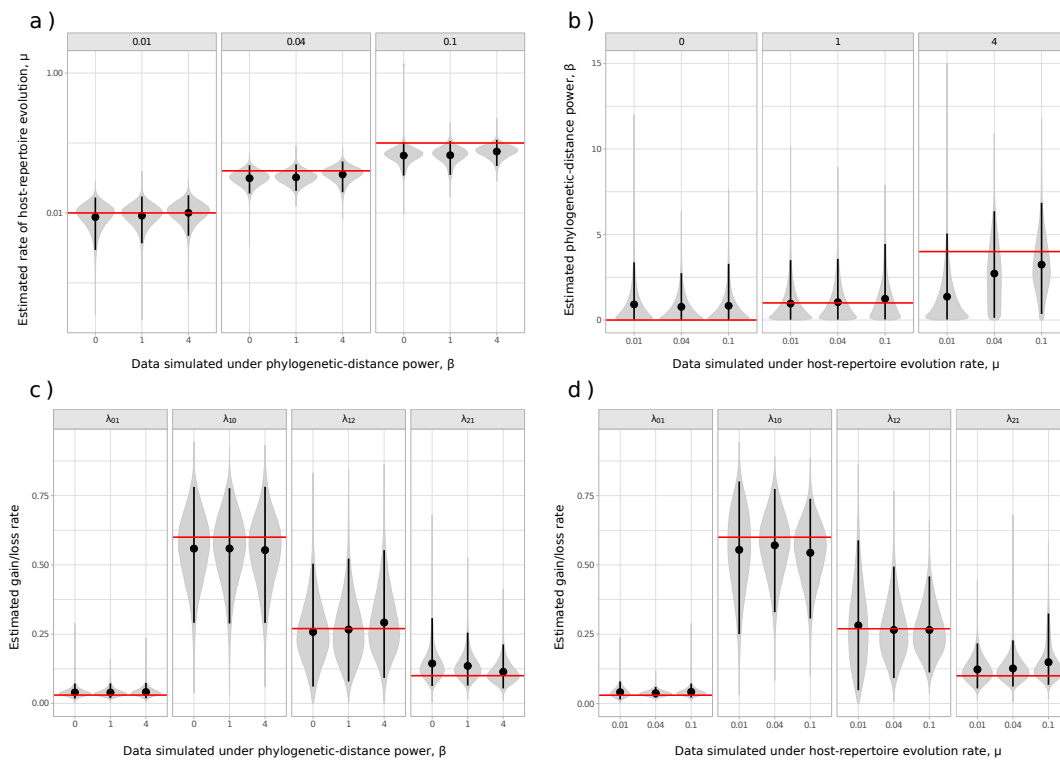
291  hosts (state 1).

Figure 3: Posterior densities of parameters in the simulation study. Panels **a** and **b** are faceted by true parameter values of $\mu$ and $\beta$, respectively. Fifty datasets were simulated for each combination of $\beta \in \{0, 1, 4\}$ and $\mu \in \{0.01, 0.04, 0.1\}$, while $\lambda_{01} = 0.03$, $\lambda_{10} = 0.6$, $\lambda_{12} = 0.27$, and $\lambda_{21} = 0.1$ were held constant. For each parameter combination, the posterior distributions of the two MCMC samples of the 50 datasets were combined. Means are represented by black dots, black vertical lines show the 95% HPD, and red horizontal lines mark the true parameter value used in the simulations. Y-axis in panel **a** is in $log_{10}$ scale for better visualization.
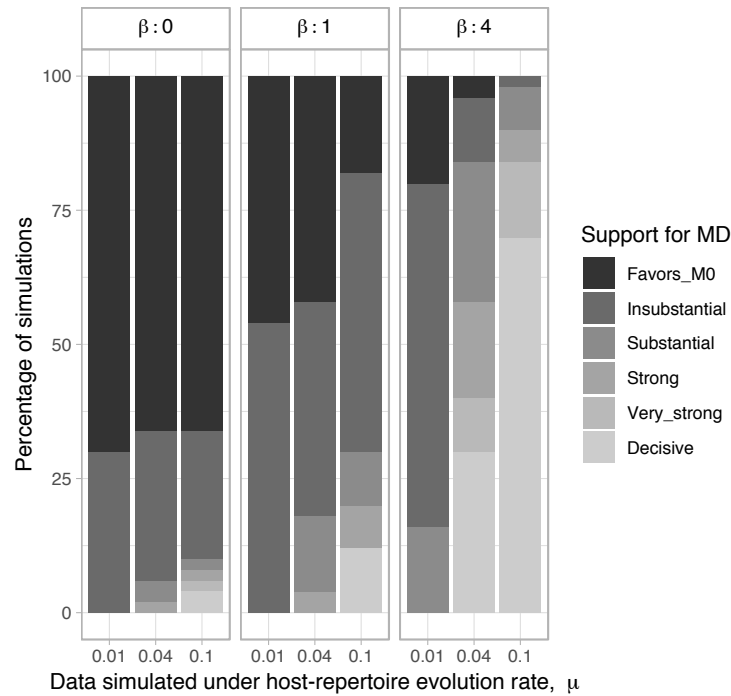
Figure 4: Distribution of Bayes factors for the simulation study. Each column corresponds to the strength of support per $2 \times 50$ MCMC analyses.
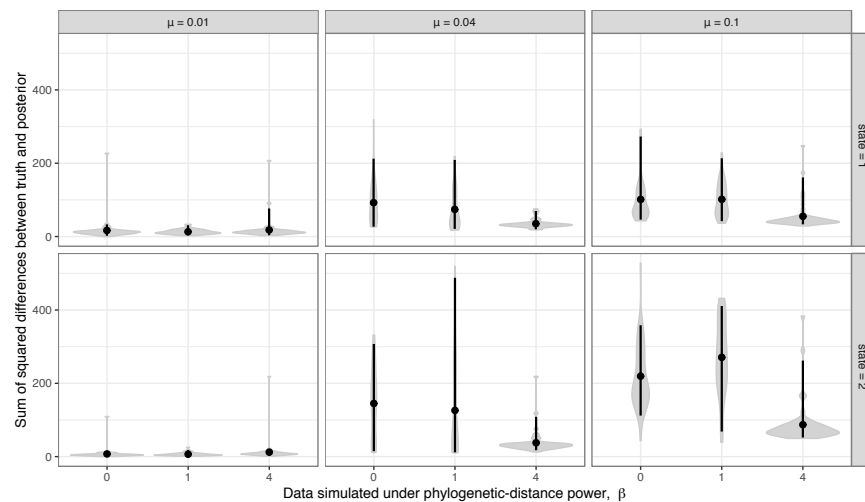


Figure 5: Errors for inferred dispersal histories of simulation study. The sum of squared differences between the posterior probability ($0 \leq P \leq 1$) and the true history ($P = 0$ or $1$) for each host and each internal node were computed per simulated dataset. Each violin plot shows the distribution of these sums for each batch of 50 simulated datasets. Means are represented by black dots, black vertical lines show the 95% CI. Values of phylogenetic-distance power ($\beta$) are shown in the x-axis, columns are separated by the host-repertoire evolution rate ($\mu$), and each row shows the error on the inference of each character state, i.e. potential host (1) or actual host (2).

292    *Empirical study.—* The estimated mean rate of host repertoire evolution for Nymphalini was

293    $\mu = 0.025$, the mean phylogenetic-distance power was $\beta = 0.51$, and the mean gain/loss rates were

294    $\lambda_{01} = 0.012$, $\lambda_{10} = 0.6$, $\lambda_{12} = 0.27$, and $\lambda_{21} = 0.12$ (Fig. 6, blue). Our method recovered similar

295    parameter estimates for the empirical dataset when omitting the intermediate state at the tips –

296    i.e. coding all potential hosts (state 1) as non-hosts (state 0): $\mu = 0.031$, $\beta = 0.39$, $\lambda_{01} = 0.001$,

297    $\lambda_{10} = 0.71$, $\lambda_{12} = 0.28$, and $\lambda_{21} = 0.01$ (Fig. 6, orange). The posterior distributions from analyses

298    with and without the intermediate state at the tips diverged the most for the rate parameters

299    associated with the transition to the intermediate state, $\lambda_{01}$ and $\lambda_{21}$. In both cases the transition

300    rate was underestimated when 1s were removed from the dataset. Bayes factors selected the

301    independence model, $M_0$, for both the full dataset (BF $= 0.43$) and when 1s were removed from
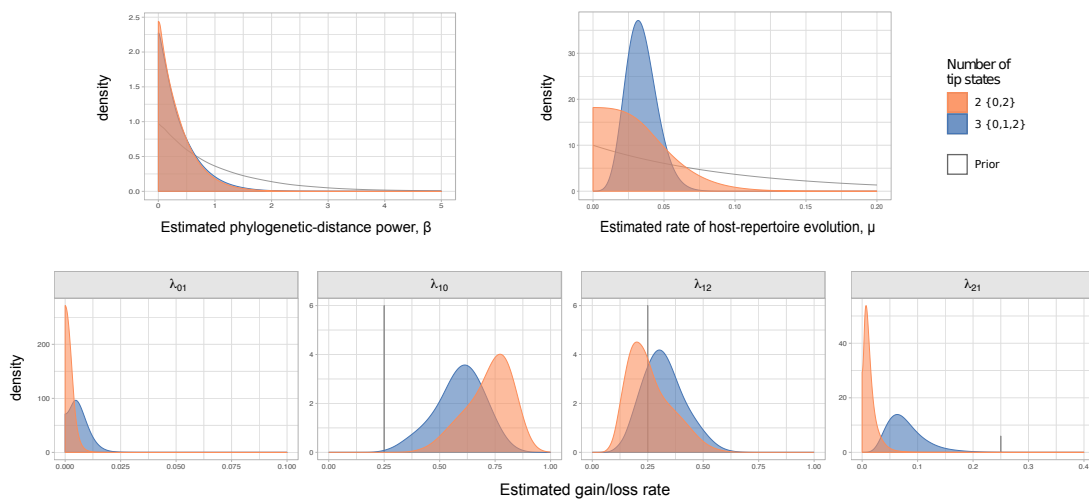
302    tip states (BF $= 0.40$).



Figure 6: Marginal posterior densities for parameters in the Nymphalini-Angiosperms study for both the full dataset (3 states at tips) and the dataset omitting the intermediate state (2 states at tips). Grey lines corresponds to the priors $\beta \sim$ Exponential(1), $\mu \sim$ Exponential(10), and $\lambda \sim$ Dirichlet(1, 1, 1, 1).

303    Finally, we reconstructed the fundamental and realized host repertoires at internal nodes

304    of the Nymphalini phylogeny based on the sampled histories during MCMC. Coevolutionary

305    histories inferred using the datasets with and without potential hosts were very similar, with

306    mean difference in interaction probability of 0.003. Thus, we only show the ancestral states

307 inferred from the full, three-state dataset (Figs. 7 and S4). To facilitate visualization of the

308 ancestral state reconstruction, we grouped the 16 parasitized host families into five modules, as

309 identified by the simulated annealing algorithm (Fig. S2). Nine families (representing three

310 modules) were inferred to be used by ancestral Nymphalini species with high probability.
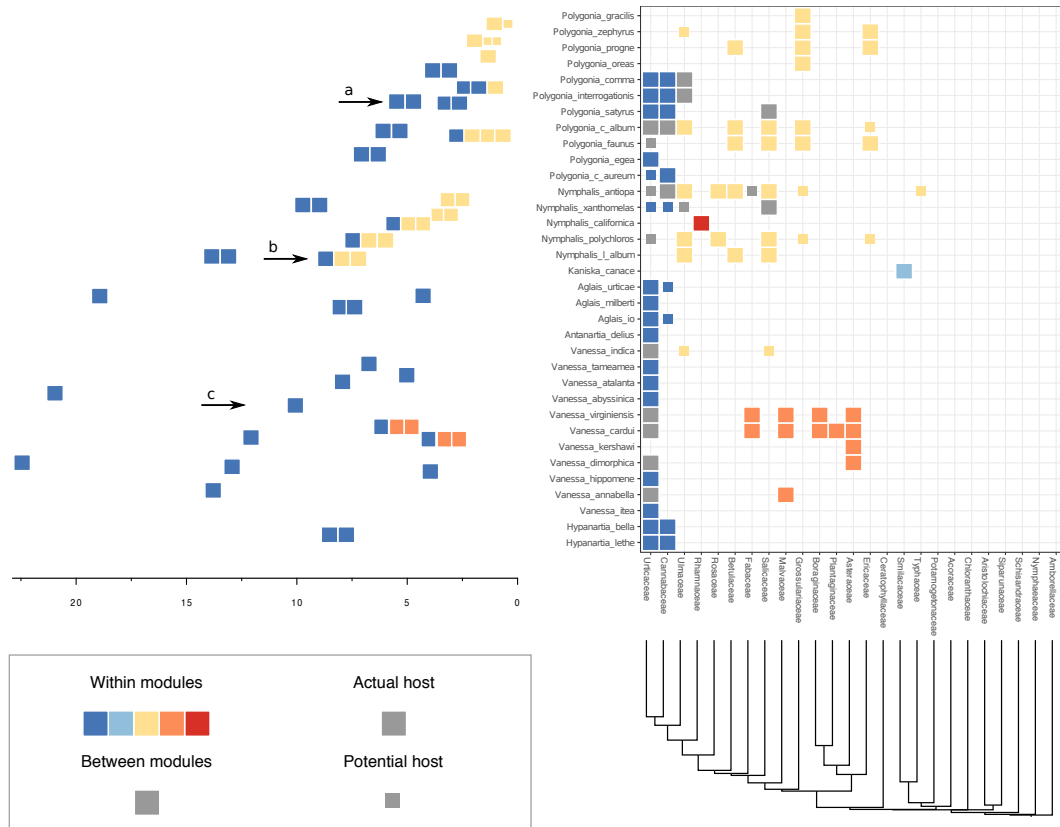


Figure 7: Evolution of butterfly-plant interactions through time. Ancestral state estimates (left) of host repertoire across the Nymphalini phylogeny are shown for interactions with more than 75% posterior probability. The x-axis under shows time before present in millions of years. Extant species interactions (right) between Nymphalini and their host plants are presented as a raster, where each square represents one interaction between a butterfly species and a host family. Colors represent different modules, i.e. groups of plants that are often hosts to the same butterflies at present time. Square size was used to differentiate between actual and potential hosts. Arrows indicate nodes shown in Fig. 8.

311 We found strong support for the association between the ancestor of all Nymphalini

312 butterflies and Urticaceae hosts (and Cannabaceae to a lesser degree, Fig. S4). All other host

313  families have been colonized in the last 15 Myr, after the divergence of the two largest clades

314  within Nymphalini, *Vanessa* and *Nymphalis + Polygonia*. Most species within *Vanessa*, both

315  extant and ancestral, are specialists on Urticaceae. *V. virginiensis* and *V. cardui* are the only

316  extant species that use more than two host families, and these hosts have likely been colonized by

317  their most recent common ancestor (node 38 in Fig. 8). On the other hand, the variation in host

318  use in the *Nymphalis + Polygonia* clade seems to be the result of host colonizations by multiple

319  species along the diversification of the clade. For example, in Fig. 8 we can see the colonization of

320  potential hosts by the ancestor of *P. c-album* and *P. faunus* (node 53) as well as strong

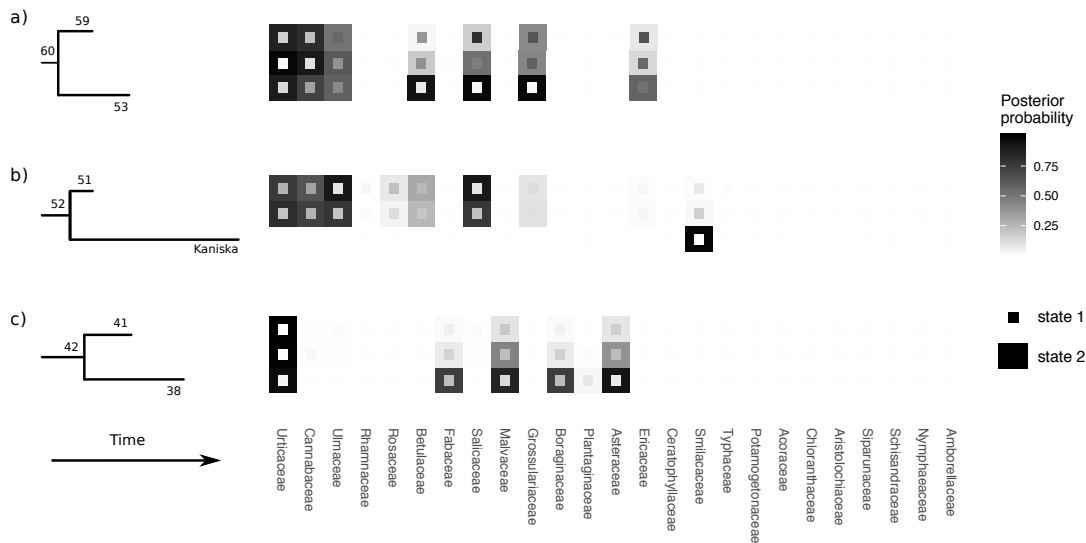321  specialization on a new host by *Kaniska canace*.



Figure 8: Host repertoires at selected nodes of the Nymphalini tree (arrows in Fig. 7). Numbers indicate the node index (compatible with Fig. S3). For the only terminal taxa depicted, *Kaniska canace*, the observed host repertoire is shown. For all other repertoires, the posterior probabilities for states 1 and 2 are shown.

## Discussion

322

323  The method we develop here to infer the evolutionary history of host-parasite associations

324  has many advantages over previous approaches. First, it is based on stochastic models and on

325  established principles of statistical inference, which means that it provides a robust framework for

326 characterizing the evolutionary processes that shape host-parasite associations and for selecting

327 among alternative coevolutionary models. Second, our model introduces the novel concept of a

328 host repertoire, which we think is an important step forward. Besides accounting for the

329 possibility of parasites having more than one host over time scales of macroevolutionary

330 significance, we can now directly infer the influence of host relatedness and host traits on the

331 process of gaining new hosts. Third, the stochastic model of host-parasite coevolution that we

332 introduce here is, to our knowledge, the first that explicitly accounts for evolution of the

333 fundamental host repertoire. By recognizing the fact that a parasite may have potential hosts in

334 addition to its actual hosts, and that the set of potential hosts may persist over time, the dynamic

335 of the model changes. What would otherwise have appeared as remarkable repeated patterns of

336 colonization of the same host lineages can now be explained as the effect of frequent transitions

337 between potential and actual hosts in an otherwise conserved host repertoire.

338      Our model can readily be extended in many interesting ways. The version we present here

339 accounts for the effect of host phylogeny by allowing the rate of host gain to depend on host

340 relatedness. For simplicity, we assumed that the number of available hosts and host relatedness

341 remain constant over geological time. This would be appropriate for a group of parasites that

342 radiated after the relevant host lineages had been formed, which is arguably the case for the

343 empirical example we chose. However, it should be relatively straightforward to extend our

344 framework to account for more complex dependencies on host phylogeny. For instance, the host

345 configurations could be modeled as changing over time, reflecting host cladogenesis.

346      Another interesting direction for future research would be to modify the particular ways in

347 which hosts and parasites coevolve. We note, for example, that Fig. 7 shows that host repertoires

348 of *Vanessa* species overlap very little with the host repertoires of *Nymphalis* + *Polygonia* species,

349 but it is not immediately clear what drives this pattern. One could design a model that allows the

350 rates of host gain and loss to be influenced by evolving host traits — like secondary metabolites,

351 growth form, or phenology, to mention a few examples relevant for insect-host plant associations

352 — in addition to relatedness among hosts. Or, one might extend the model to allow closely

353 related parasite lineages to competitively exclude one another from host usage, similar to how

354 competing lineages might exclude one another from geographical regions (Quintero and Landis

355 2019). Finally, one might introduce a biogeographical component to the coevolutionary process,

356 requiring parasites to be in sympatry with their actual hosts, while allowing parasites to be in

357 sympatry or allopatry with their potential hosts. Statistically comparing such model variants will

358 help illuminate drivers of host-parasite co-evolution.

359 A potential concern with our approach is that already the basic version of the model is

360 fairly parameter-rich. Given the type and amount of data that we can likely collect on

361 host-parasite associations, is there enough statistical power to select among the models of interest?

362 And is it possible to infer the model parameters of interest with a reasonable degree of accuracy?

363 Overall, our results are encouraging in this respect. The simulations indicate that it is

364 possible to infer the true parameter values of the basic model regardless of the level of

365 phylogenetic conservatism in both parasites and hosts (Fig. 3). When the rates of colonization of

366 new hosts are strongly dependent on the phylogenetic relatedness of hosts, then we are also able

367 to distinguish between models with or without host relatedness effects using Bayes factors (Fig.

368 4). However, our ability to select the correct model decreases when the effect of host phylogenetic

369 relatedness is low ($\beta \leq 1$), that is, when models become more similar. Further studies will have to

370 show to what extent the sensitivity of the model test can be increased by selecting appropriate

371 priors and improving the sampling of parameter space close to the boundary condition satisfying

372 the restricted model. One option is to relax the assumption that $\beta$ is non-negative, which would

373 simplify the sampling of values close to $\beta = 0$. It will also be important to explore how dataset

374 sizes and tree shapes, for both hosts and parasites, influence our ability to distinguish the models

375 when the effect of host phylogeny is small.

376 Importantly, the empirical analysis indicates that the method is able to model the

377 evolution of fundamental and realized host repertoires even when the information about potential

378 hosts is lacking. This significantly increases the applicability of our method, as information about

379 fundamental host repertoires is missing for most host-parasite systems. Potential host data is

380 difficult to collect, as it requires experimental testing of a large number of potential host-parasite

381 pairs. A possible improvement of our method, which we did not explore here, would be to model

382 uncertainty in the observations of non-hosts when data on potential hosts are missing. That is, if

383 we had no information about a host species being used by a particular parasite, we would

384 translate that to a certain probability $p$ of the species actually being a non-host, and a

385 complementary probability $1 - p$ of it being a potential host (Kuhner and McGill 2014). Modeling

386 this observational uncertainty could help reduce the bias in parameter estimates that we observed

387 when data on potential hosts were missing and all 0 states in the dataset were inappropriately

388 treated as true non-hosts. This extension would also allow us to make predictions about host use

389 abilities in extant parasites. These predictions could then inform experiments that aim to

390 characterize fundamental host repertoires.

391 We demonstrated the empirical application of our approach with a Bayesian inference of

392 the coevolutionary history between 34 Nymphalini butterflies and 25 angiosperm families. We

393 estimated the rate of host repertoire evolution along each branch of the butterfly tree as being

394 between 0.33 and 0.93 events per million years. Bayes factors favored the independence model,

395 where the probability of gaining a given hosts is not affected by the phylogenetic distance between

396 hosts. As explained above, this does not necessarily mean that host relatedness plays no role, only

397 that the effect is not large enough for us to detect it with the current approach and the given data.

398 Estimates of gain and loss rates were not symmetric, and the rates also varied between

399 states. According to our results, gain of the ability to use a host, $\lambda_{01}$, is very rare (0.5% to 1.9%

400 of overall rate), whereas loss is common (47% to 73% of overall rate). On the other hand,

401 transition rates between states 1 and 2 were more symmetric and gain is more common than loss

402 ($\lambda_{12}$ between 15% and 39%; $\lambda_{21}$ between 6% and 18% of overall rate). These rate estimates

403 support the idea that the use of the same host lineage by multiple, phylogenetically widespread

404 butterfly lineages is more likely explained by recolonization of hosts that have been used in the

405 past (recurrence homoplasy), that is, by transitions between actual and potential hosts, rather

406 than by completely independent colonizations of the same host (Janz et al. 2001). Note that

407 alternative scenarios that have been proposed in the literature to explain the evolution of

408 Nymphalini host plant preferences, for instance by involving narrow ancestral host plant ranges

409 and repeated independent colonization events, are also allowed by our model, but they are

410 inferred to be much less likely than the conservative host repertoire scenario. Yet, because the

411 potential host state is exited at the highest rate, the rate estimates also suggest that parasites do

412 not retain their potential host relationships for prolonged periods of time. The moderate rates of

413 transitions between potential and actual host states and the high departure rate from the

414 potential host state together help explain why phylogenetic "pulses" of recurrent host acquisition

415 manifest in some lineages but not others.

416      For example, the use of Grossulariaceae by two non-sister clades within *Polygonia* is best

417 explained by a scenario where Grossulariaceae was a potential host for the ancestral species (node

418 60 in Fig. 8) and was subsequently gained as an actual host twice (at nodes 53 and 58, Fig. S4).

419 The ability to use Salicaceae host plants seems to be even older. Salicaceae was likely a potential

420 host for the ancestor of *Nymphalis + Polygonia* and later became an actual host in three different

421 parts of the clade. If potential hosts were not explicitly modeled here, these transitions would

422 look like three independent colonizations of a plant group that is very distant from the ancestral

423 host (Salicaceae and Urticaceae diverged about 90 Ma). Instead, we could show that what might

424 appear as big and sudden host shifts, are in fact the result of retention of ancestral host use

425 abilities.

426      Understanding how ecological interactions change is crucial if we want to predict both

427 short and long-term consequences of global mixing of biota (Hoberg and Brooks 2015).

428 Host-parasite interactions are of particular interest given the risk of emerging diseases, which can

429 affect human populations directly and indirectly through their effects on crop species and wildlife

430 (Brooks et al. 2014). Our method was designed to quantify changes in host-parasite associations

431 by modeling the process of gaining and losing hosts, thus allowing us to make predictions based

432 on host-parasite history. Hopefully, our approach will not only generate deeper insights into the

433 evolutionary dynamics of host-parasite associations but also help humankind mitigate some of the

434    risks incurred by current environmental change.

## FUNDING

439                                                            *

440    Acha, P. N. and B. Szyfres. 2003. Zoonoses and communicable diseases common to man and

441        animals vol. 580. Pan American Health Org.

442    Barber, M. J. 2007. Modularity and community detection in bipartite networks. Physical review.

443        E, Statistical, nonlinear, and soft matter physics 76:066102.

444    Braga, M. P., P. R. Guimarães Jr, C. W. Wheat, S. Nylin, and N. Janz. 2018. Unifying

445        host-associated diversification processes using butterfly-plant networks. Nature communications

446        9.

447    Brooks, D. R. 1979. Testing the Context and Extent of Host-Parasite Coevolution. Systematic

448        Biology 28:299–307.

449    Brooks, D. R., E. P. Hoberg, W. A. Boeger, S. L. Gardner, K. E. Galbreath, D. Herczeg, H. H.

450        Mejia-Madrid, S. E. Racz, and A. T. Dursahinhan. 2014. Finding Them Before They Find Us:

451        Informatics, Parasites, and Environments in Accelerating Climate Change. Comparative

452        Parasitology 81:155–164.

453    Calatayud, J., J. L. Hórreo, J. Madrigal-González, A. Migeon, M. Á. Rodríguez, S. Magalhães,

454        and J. Hortal. 2016. Geography and major host evolutionary transitions shape the resource use

455        of plant parasites. Proceedings of the National Academy of Sciences 113:201608381–9845.

de Vienne, D. M., G. Refregier, M. Lopez-Villavicencio, A. Tellier, M. E. Hood, and T. Giraud. 2013. Cospeciation vs host-shift speciation: methods for testing, evidence from natural associations and relation to coevolution. New Phytologist 198:347–385.

Elton, C. 1946. Competition and the Structure of Ecological Communities. Journal of Animal Ecology 15:54–68.

Fisher, M. C., T. W. Garner, and S. F. Walker. 2009. Global emergence of *Batrachochytrium dendrobatidis* and amphibian chytridiomycosis in space, time, and host. Annual Review of Microbiology 63:291–310.

Forister, M. L., V. Novotny, A. K. Panorska, L. Baje, Y. Basset, P. T. Butterill, L. Cizek, P. D. Coley, F. Dem, I. R. Diniz, P. Drozd, M. Fox, A. E. Glassmire, R. Hazen, J. Hrcek, J. P. Jahner, O. Kaman, T. J. Kozubowski, T. A. Kursar, O. T. Lewis, J. Lill, R. J. Marquis, S. E. Miller, H. C. Morais, M. Murakami, H. Nickel, N. A. Pardikes, R. E. Ricklefs, M. S. Singer, A. M. Smilanich, J. O. Stireman, S. Villamarín-Cortez, S. Vodka, M. Volf, D. L. Wagner, T. Walla, G. D. Weiblen, and L. A. Dyer. 2015. The global distribution of diet breadth in insect herbivores. Proceedings of the National Academy of Sciences 112:442–447.

Gelman, A. and D. B. Rubin. 1992. Inference from Iterative Simulation Using Multiple Sequences. Statistical Science 7:457–472.

Hahn, B. H., G. M. Shaw, K. M. De Cock, and P. M. Sharp. 2000. Aids as a zoonosis: scientific and public health implications. Science 287:607–614.

Hardy, N. B. 2017. Do plant-eating insect lineages pass through phases of host-use generalism during speciation and host switching? Phylogenetic evidence. Evolution 71:2100–2109.

Hastings, W. K. 1970. Monte Carlo sampling methods using Markov chains and their applications. Biometrika 57:97–109.

Hoberg, E. P. and D. R. Brooks. 2015. Evolution in action: climate change, biodiversity dynamics

and emerging infectious disease. Philosophical transactions of the Royal Society of London. Series B, Biological sciences 370:–20130553.

Höhna, S., M. J. Landis, T. A. Heath, B. Boussau, N. Lartillot, B. R. Moore, J. P. Huelsenbeck, and F. Ronquist. 2016. RevBayes: Bayesian Phylogenetic Inference Using Graphical Models and an Interactive Model-Specification Language. Systematic Biology 65:726–736.

Huelsenbeck, J. P., B. Rannala, and Z. H. Yang. 1997. Statistical tests of host-parasite cospeciation. Evolution 51:410–419.

Janz, N., M. P. Braga, N. Wahlberg, and S. Nylin. 2016. On oscillations and flutterings—A reply to Hamm and Fordyce. Evolution 70:1150–1155.

Janz, N., K. Nyblom, and S. Nylin. 2001. Evolutionary dynamics of host-plant specialization: a case study of the tribe Nymphalini. Evolution 55:783–796.

Janz, N. and S. Nylin. 2008. The oscillation hypothesis of host-plant range and speciation. Pages 203–215 *in* Specialization, speciation, and radiation: the evolutionary biology of herbivorous insects (K. J. Tilmon, ed.). California Univ. Press, California.

Jeffreys, H. 1961. The Theory of Probability. OUP Oxford.

Klassen, G. J. 1992. Coevolution: A History of the Macroevolutionary Approach to Studying Host-Parasite Associations. The Journal of Parasitology 78:573.

Kuhner, M. K. and J. McGill. 2014. Correcting for sequencing error in maximum likelihood phylogeny inference. G3: Genes, Genomes, Genetics 4:2545–2552.

Landis, M. J., N. J. Matzke, B. R. Moore, and J. P. Huelsenbeck. 2013. Bayesian analysis of biogeography when the number of areas is large. Systematic Biology 62:789–804.

Larose, C., S. Rasmann, and T. Schwander. 2019. Evolutionary dynamics of specialisation in herbivorous stick insects. Ecology Letters 22:354–364.

503 Legendre, P., Y. Desdevises, and E. Bazin. 2002. A Statistical Test for Host–Parasite Coevolution.
504     Systematic Biology 51:217–234.

505 Magallón, S., S. Gómez-Acevedo, L. L. Sánchez-Reyes, and T. Hernández-Hernández. 2015. A
506     metacalibrated time-tree documents the early rise of flowering plant phylogenetic diversity.
507     New Phytologist 207:437–453.

508 Marquitti, F. M. D., P. R. Guimarães, M. M. Pires, and L. F. Bittencourt. 2014. MODULAR:
509     software for the autonomous computation of modularity in large network sets. Ecography
510     37:221–224.

511 Newman, M. E. J. and M. Girvan. 2004. Finding and evaluating community structure in
512     networks. Physical review. E, Statistical, nonlinear, and soft matter physics 69:026113.

513 Nielsen, R. 2002. Mapping Mutations on Phylogenies. Systematic Biology 51:729–739.

514 Nosil, P. 2002. Transition rates between specialization and generalization in phytophagous insects.
515     Evolution 56:1701–1706.

516 Nylin, S., S. Agosta, S. Bensch, W. A. Boeger, M. P. Braga, D. R. Brooks, M. L. Forister, P. A.
517     Hambäck, E. P. Hoberg, T. Nyman, A. Schäpers, A. L. Stigall, C. W. Wheat, M. Österling, and
518     N. Janz. 2018. Embracing Colonizations: A New Paradigm for Species Association Dynamics.
519     Trends in Ecology & Evolution 33:4–14.

520 Nylin, S., J. Slove, and N. Janz. 2014. Host plant utilization, host range oscillations and
521     diversification in nymphalid butterflies: a phylogenetic investigation. Evolution 68:105–124.

522 Plummer, M., N. Best, K. Cowles, K. Vines, and 2006. 2006. CODA: convergence diagnosis and
523     output analysis for MCMC. R News 6:7–11.

524 Quintero, I. and M. J. Landis. 2019. Interdependent Phenotypic and Biogeographic Evolution
525     Driven by Biotic Interactions. BioRxiv Page 560912.

526 Ree, R. H., B. R. Moore, C. O. Webb, and M. J. Donoghue. 2005. A likelihood framework for

527     inferring the evolution of geographic range on phylogenetic trees. Evolution 59:2299–2311.

528 Ree, R. H. and S. A. Smith. 2008. Maximum likelihood inference of geographic range evolution by

529     dispersal, local extinction, and cladogenesis. Systematic Biology 57:4–14.

530 Robinson, D. M. 2003. Protein Evolution with Dependence Among Codons Due to Tertiary

531     Structure. Molecular Biology and Evolution 20:1692–1704.

532 Ronquist, F. 2003. Parsimony analysis of coevolving species associations . Pages 22–64 *in* Tangled

533     trees: Phylogeny, cospeciation and coevolution. (R. D. M. Page, ed.). University of Chicago

534     Press, Chicago.

535 Subbarao, K., A. Klimov, J. Katz, H. Regnery, W. Lim, H. Hall, M. Perdue, D. Swayne,

536     C. Bender, J. Huang, et al. 1998. Characterization of an avian influenza A (H5N1) virus

537     isolated from a child with a fatal respiratory illness. Science 279:393–396.

538 Suchard, M. A., R. E. Weiss, and J. S. Sinsheimer. 2001. Bayesian selection of continuous-time

539     Markov chain evolutionary models. Molecular Biology and Evolution 18:1001–1013.

540 Verdinelli, I. and L. Wasserman. 1995. Computing Bayes Factors Using a Generalization of the

541     Savage-Dickey Density Ratio. Journal of the American Statistical Association 90:614–618.