

1 Complex ecological phenotypes on phylogenetic trees: a hidden Markov model for comparative
2 analysis of multivariate count data

3

4 Michael C. Grundler^{1,*} and Daniel L. Rabosky¹

5

6 *¹Museum of Zoology and Department of Ecology and Evolutionary Biology, University of*

7 *Michigan, Ann Arbor, MI USA 48109*

8

9 **Corresponding author: mgru@umich.edu*

10

11

12

13

14

15

16

17

18

19

20

21

22

23

GRUNDLER AND RABOSKY

24 ABSTRACT

25 The evolutionary dynamics of complex ecological traits – including multistate
26 representations of diet, habitat, and behavior – remain poorly understood. Reconstructing the
27 tempo, mode, and historical sequence of transitions involving such traits poses many challenges
28 for comparative biologists, owing to their multidimensional nature and intraspecific variability.
29 Continuous-time Markov chains (CTMC) are commonly used to model ecological niche
30 evolution on phylogenetic trees but are limited by the assumption that taxa are monomorphic and
31 that states are univariate categorical variables. Thus, a necessary first step when using standard
32 CTMC models is to categorize species into a pre-determined number of ecological states. This
33 approach potentially confounds interpretation of state assignments with effects of sampling
34 variation because it does not directly incorporate empirical observations of resource use into the
35 statistical inference model. The neglect of sampling variation, along with univariate
36 representations of true multivariate phenotypes, potentially leads to the distortion and loss of
37 information, with substantial implications for downstream macroevolutionary analyses. In this
38 study, we develop a hidden Markov model using a Dirichlet-multinomial framework to model
39 resource use evolution on phylogenetic trees. Unlike existing CTMC implementations, states are
40 unobserved probability distributions from which observed data are sampled. Our approach is
41 expressly designed to model ecological traits that are intra-specifically variable and to account
42 for uncertainty in state assignments of terminal taxa arising from effects of sampling variation.
43 The method uses multivariate count data for individual species to simultaneously infer the
44 number of ecological states, the proportional utilization of different resources by different states,
45 and the phylogenetic distribution of ecological states among living species and their ancestors.

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

46 The method is general and may be applied to any data expressible as a set of observational
47 counts from different categories.

48

49 Keywords: Ecological niche evolution; intraspecific variation; hidden Markov model;
50 macroevolution; comparative methods; Dirichlet-multinomial

51

52 Most species in the natural world make use of multiple, categorically-distinct types of
53 ecological resources. Many butterfly species use multiple host plants, for example (Ehrlich &
54 Raven 1964; Robinson 1999). Insectivorous warblers in temperate North America use multiple
55 distinct microhabitats and foraging behaviors (MacArthur 1958), as do honeyeaters in mesic and
56 arid Australia (Miller et al. 2017). The evolution of novel patterns of resource use can impact
57 phenotypic evolution (Martin & Wainwright 2011; Davis et al. 2016), diversification (Mitter et
58 al. 1988; Givnish et al. 2014), community assembly (Losos et al. 2003; Gillespie 2004), and
59 ecosystem function (Harmon et al. 2009; Bassar et al. 2010). Consequently, there has been
60 substantial interest in understanding how ecological traits related to resource use evolve and for
61 exploring their impacts on other evolutionary and ecological phenomena (Vrba 1987; Futuyma &
62 Moreno 1988; Forister et al. 2012; Price et al. 2012; Burin et al. 2016).

63 Making inferences about the evolutionary dynamics of resource use, however, first
64 requires summarizing the complex patterns of variation observed among taxa into traits that can
65 be modeled on phylogenetic trees. It is widely recognized that the real-world complexities of
66 resource use are not adequately described by a set of categorical variables (Hardy & Linder
67 2005; Hardy 2006). Nonetheless, it is also true that major differences in resource use can
68 sometimes be summed up in a small set of ecological states, a point made by Mitter et al. (1988)

GRUNDLER AND RABOSKY

69 in their study of phytophagy and insect diversification. For this reason, continuous-time Markov
70 chain (CTMC) models, which require classifying species into a set of character states, have
71 become commonplace in macroevolutionary studies of ecological trait evolution (Kelley &
72 Farrell 1998; Nosil 2002; Price et al. 2012; Hardy & Otto 2014; Cantalapiedra et al. 2014; Burin
73 et al. 2016). CTMC models describe a stochastic process for evolutionary transitions among a set
74 of character states and are used to infer ancestral states and evolutionary rates, and to perform
75 model-based hypothesis tests (O’Meara 2012).

76 The utility of continuous-time Markov chains for studying the evolutionary dynamics of
77 resource use is limited by the modeling assumption that taxa are monomorphic for ecological
78 states (Hardy & Linder 2005; Hardy 2006). As a practical solution, most empirical studies define
79 one or more generalized states to accommodate species that use multiple resource types and that
80 therefore cannot be characterized as specialists for a particular resource (Alencar et al. 2013;
81 Price et al. 2012; Burin et al. 2016; Gajdzik et al. 2019). Another solution, rather than classifying
82 each species as a specialist or a generalist, represents each resource category with a binary score
83 of present or absent (Janz et al. 2001; Colston et al. 2010; Hardy 2017). In this case, the
84 ecological state of a species is the set resources scored as present. Each of these approaches is
85 one solution to the modeling challenge posed by intraspecific variation in resource use, but both
86 solutions neglect variation in the relative importance of different resources for different taxa.
87 Consequently, species classified in single state can nonetheless exhibit substantial differences in
88 patterns of resource use, creating challenges for interpreting evolutionary transitions among
89 character states as well as for understanding links between character state evolution and
90 diversification.

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

91 Another limitation of continuous-time Markov chains for modeling resource use
92 evolution emerges from the fact that species are classified into ecological states without regard
93 for the quality and quantity of information available to perform the classification exercise. As an
94 example, species with few ecological observations might be classified as specialists for a
95 particular resource, when their apparent specialization is strictly a function of the small number
96 of ecological observations available for the taxon. More generally, by failing to use a statistical
97 model for making resource state assignments, we neglect a major source of uncertainty in our
98 data: the uneven and incomplete knowledge of resource use across different taxa. This
99 uncertainty, in turn, has substantial implications for how we project patterns of resource use onto
100 a set of resource states. By failing to account for uneven and finite sample sizes characteristic of
101 empirical data on resource use we cannot be certain if state assignments reflect true similarities
102 or differences in resource use or are merely the expected outcome of sampling variation.

103 Consider the simple four-species example in Figure 1. Panels (i) and (ii) illustrate the true
104 resource states and their phylogenetic distribution across a set of four species and their ancestors.
105 Here, an ancestral specialist evolved a generalist diet via a single transition (panel ii), such that
106 there are two extant species with the ancestral specialist diet (species X and Y) and two with the
107 derived generalist diet (species P and Q). In panels (iii) and (iv) the relative importance estimates
108 of three food resource categories in the diets of four species are used to classify each species into
109 one of two diet states. These relative importance estimates are based on uneven, and in some
110 cases quite small, sample sizes, consistent with many empirical datasets (Vitt & Vangilder 1983;
111 Shine 1994; Alencar et al. 2013). In panel (v) we imagine repeating the state assignment process
112 on independent datasets while holding the samples sizes fixed to those in panel (iii), which
113 reveals that both the initial state assignments and the number of states from (iv) are highly

GRUNDLER AND RABOSKY

114 sensitive to real-world levels of sampling variation. This has obvious implications for
115 downstream macroevolutionary analyses. There is a serious risk of conflating different (similar)
116 state assignments with different (similar) diets when the differences (similarities) are expected,
117 even in the absence of true differences (similarities), from sampling variation alone. In the
118 analyses that underlie Figure 1, we find that more than 70 percent of tip state classifications do
119 not match the true pattern of resource use.

120 Is this a problem in practice? This issue is difficult to assess because few studies provide
121 information about the sample sizes that underlie state assignments. In most cases, ecological
122 states are simply asserted as known. It is also important to emphasize that the specific problem in
123 Figure 1 is an outcome of a more general problem: standard CTMC models have a limited ability
124 to model complex ecological phenotypes because of the assumption that states in the model are
125 categorical variables. While it is true that CTMC models work with a countable state space, it is
126 not true that the states of the system must represent categorical variables. In a hidden Markov
127 model, the observed data are assumed to be the outcome of a CTMC where the states are not
128 directly observable. Instead, states are probability distributions from which observed data are
129 sampled. Although hidden state CTMC models already used in macroevolution (Marazzi et al.
130 2012; Beaulieu et al. 2013; Beaulieu & O'Meara 2016; Caetano et al. 2018) can be interpreted
131 like this they generally are not because in these cases the observed data are categorical variables
132 and the hidden states are indistinguishable from the observed categories. Instead, hidden states
133 are interpreted as unobserved factors that affect rates of change or rates of diversification.
134 Because of this, the potential flexibility of hidden state models for modeling complex phenotypes
135 remains poorly explored.

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

136 In this paper we use the formulation of hidden states as probability distributions to
137 develop a CTMC model for studying the evolutionary dynamics of ecological resource use on
138 phylogenetic trees. Our approach is explicitly designed to model resource traits that are intra-
139 specifically variable and to account for uncertainty in ecological state assignments of terminal
140 taxa arising from effects of sampling variation. We assume that each state is an unobserved
141 (latent) multinomial distribution and that observed data are sampled outcomes from these latent
142 distributions (see panels (i) through (iii) of Fig. 1). The number of states in the model and the
143 states themselves are not directly observed and are estimated from the data. Using simulations
144 and an empirical dataset of snake diets, we show how the method can use observational counts to
145 simultaneously infer the number of resource states, the proportional utilization of resources by
146 different states, and the phylogenetic distribution of ecological states among living species and
147 their ancestors. The method is general and applicable to any data expressible as a set of
148 observational counts from different resource categories.

149

150 MATERIALS & METHODS

151 *Model description*

152 We assume that the data for each species are represented by a vector of J category counts.
153 Each category is a resource (e.g. a diet or habitat component), and each count represents the
154 number of observations of a species utilizing a particular resource. Each node in a phylogeny is
155 to be placed into one of K distinct resource states. States are unobserved, even at the tips of a
156 phylogeny: the observed data consist of sampled outcomes from an underlying (latent)
157 multinomial distribution that represents the state for each species. All count data are drawn
158 independently from their respective states and counts from each state are also independent of one

GRUNDLER AND RABOSKY

159 another. We assume that the multinomial parameters for each state are drawn from a common
160 Dirichlet distribution with parameter β . This parameterization allows us to analytically
161 marginalize over the unknown multinomial parameters underlying each state so that the
162 likelihood of the observed data is the product of K independent Dirichlet-multinomial
163 distributions (Appendix). The parameter β acts as a vector of pseudo-counts. Higher values
164 require more data for the model to discriminate two samples as having originated from different
165 states. Letting X denote the resource state assignments for nodes in the phylogeny, the likelihood
166 of the set of count data D_k generated from state k is,

$$167 \quad p(D_k|X, \beta) = \frac{\Gamma(J\beta)}{\Gamma(n_k + J\beta)} \frac{\prod_j \Gamma(n_k^j + \beta)}{\Gamma(\beta)^J} \quad (1)$$

168 where n_k is the total number of observations generated from state k and n_k^j is the subset
169 of those observations that represent utilization of resource category j . The full likelihood for the
170 count data is just,

$$171 \quad p(D|X, \beta) = \prod_k p(D_k|X, \beta) \quad (2)$$

172 This model for count data is closely related to topic models of word composition in a
173 collection of text documents (Blei et al. 2003; Yin and Wang 2014) and to population genetic
174 models of allele frequency composition in a set of populations (e.g., program STRUCTURE:
175 Pritchard et al. 2000). The key difference here is that the state assigned to a taxon is the outcome
176 of evolution and is not independent of the states of other lineages. Conceptually this is similar to
177 phylogenetic threshold models, where the full likelihood combines a probability model for the
178 evolution of an unobserved variable and a probability model for sampling the observed data
179 conditioned on the set of unobserved variables (Felsenstein 2012; Revell 2014). We model
180 evolution as a Poisson process where the rate of change is the same between all states (i.e. there

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

181 is no evolutionary trend in the model) but varies among lineages. We introduce two mechanisms
182 for accommodating this rate variation.

183 The first mechanism takes advantage of the random local clocks model introduced by
184 Drummond and Suchard (2010). In this framework, there is an overall rate of evolution Λ and an
185 unknown number of lineages that deviate from this rate by a set of multiplicative constants.
186 Specifically, the root node is defined to have a relative rate of 1 while the relative rates of all
187 other nodes are equal to the relative rate of their ancestor multiplied by a branch-specific positive
188 rate multiplier. Complexity is controlled through the use of a prior that makes it unlikely for
189 many of these multipliers to differ from unity. Under the fully symmetric Poisson model with
190 random local clocks the probability of change across an ancestral-descendant branch is,

$$191 \quad p(X_i | X_{pa(i)}, \Lambda, r_i) = (1 - e^{-K\Lambda r_i t_i}) \frac{1}{K} + \delta_{X_i X_{pa(i)}} e^{-K\Lambda r_i t_i} \quad (3)$$

192 where X_i is the state of node i , $X_{pa(i)}$ is the state of i 's ancestor, r_i is the branch-specific
193 normalized relative rate of evolution, t_i is the length of the branch in units of time, and $\delta_{X_i X_{pa(i)}}$
194 is the Kronecker delta. For simplicity, we may occasionally notate transition probabilities by p_{ii}
195 and p_{ij} to indicate the probability that a descendant's state is the same as, or different than, the
196 state of its ancestor. The likelihood of the node states is just,

$$197 \quad p(X | \Lambda, r) = \prod_i p(X_i | X_{pa(i)}, \Lambda, r_i) \quad (4)$$

198 where the product is taken over all nodes and $p(X_i | X_{pa(i)}, \Lambda, r_i) \equiv \frac{1}{K}$ if node i is the root.

199 The normalized rates effectively expand and contract the temporal durations of the branches.
200 They are derived by scaling the relative rates in such a way that the total absolute time in which
201 evolution has had to occur is held constant even as the effective time is allowed to vary over
202 phylogeny (that is, $\sum_i t_i = \sum_i r_i t_i$).

GRUNDLER AND RABOSKY

203 The second mechanism for accommodating rate heterogeneity is essentially a saturated
204 version of the random local clocks model where each branch has a unique rate of evolution.
205 Following Huelsenbeck et al. (2008), this allows us to model branch-specific rates as nuisance
206 parameters drawn independently from a Gamma distribution with parameter vector $(\alpha, 1)$. This
207 model induces the same distribution of node states as a model where the number of expected
208 character state changes along a branch is the same for all branches (Appendix). This has
209 elsewhere been termed the ultra-common mechanism model (Steel 2011) to mark its contrast
210 with the no-common mechanism model (Tuffley and Steel 1997) from which it derives. In this
211 case the probability of change across an ancestral-descendant branch is,

$$212 \quad p(X_i | X_{pa(i)}, \alpha) = \left(1 - \frac{1}{(K/(K-1) + 1)^\alpha}\right) \frac{1}{K} + \delta_{X_i X_{pa(i)}} \frac{1}{(K/(K-1) + 1)^\alpha} \quad (5)$$

213 Phylogenetic signal is controlled by the parameter α , which is equal to the expected number of
214 state changes that occur from ancestor to descendant. As $\alpha \rightarrow 0$, phylogenetic signal approaches
215 1 because descendants almost surely resemble their ancestors. As $\alpha \rightarrow \infty$, phylogenetic signal
216 approaches 0 because a descendant's state becomes independent of its ancestor's state and
217 resembles a random draw from a discrete uniform distribution. The likelihood of the node states
218 is just,

$$219 \quad p(X|\alpha) = \frac{1}{K} p_{ii}^n p_{ij}^m \quad (6)$$

220 where n is the number of nodes with the same state as their ancestor, m is the number of
221 nodes with a different state than their ancestor, and the factor $\frac{1}{K}$ accounts for the root state
222 probability.

223

224 *Bayesian inference*

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

225 We simulated the posterior distribution of node states and model parameters using the
 226 Metropolis-Hastings algorithm (Hastings 1970). The different proposal mechanisms are
 227 described below.

228 *Updating node states.*— A Gibbs update mechanism is used for proposing changes to
 229 node states. The full conditional distribution for the state of a node can be written as,

$$\begin{aligned}
 230 \quad p(X_i|X_{-i}, D, \beta, \theta) &= \frac{p(D, X|\beta, \theta)}{p(D, X_{-i}|\beta, \theta)} \\
 231 &\propto \frac{p(D, X|\beta, \theta)}{p(D_{-i}, X_{-i}|\beta, \theta)} \\
 232 &= \frac{p(D|X, \beta)}{p(D_{-i}|X_{-i}, \beta)} \frac{p(X|\theta)}{p(X_{-i}|\theta)} \quad (7)
 \end{aligned}$$

233 where the symbol $-i$ denotes the exclusion of node i . Here, depending on whether (4) or
 234 (6) is used as the likelihood model for node states, θ is equal to (Λ, r) or α , respectively. Because
 235 changing a node state only affects the branches incident to the affected node, all terms in the
 236 ratio $\frac{p(X|\theta)}{p(X_{-i}|\theta)}$ not involving those branches cancel and it simplifies to,

$$237 \quad \frac{p(X|\theta)}{p(X_{-i}|\theta)} = p(X_i|X_{pa(i)}, \theta) \prod_k p(X_{ch(i,k)}|X_i, \theta)$$

238 where $X_{ch(i,k)}$ is the state of the k -th immediate descendant of node i . Similarly, altering
 239 the state of a node only affects the likelihood of count data associated with a single state so that
 240 factors in the ratio $\frac{p(D|\beta)}{p(D_{-i}|\beta)}$ not involving, say, state k cancel, and it simplifies to,

$$\begin{aligned}
 241 \quad \frac{p(D|X, \beta)}{p(D_{-i}|X_{-i}, \beta)} &= \frac{p(D_k|X, \beta)}{p(D_{k-i}|X_{-i}, \beta)} \\
 242 &= \frac{\prod_j \Gamma(n_{k-i}^j + N_i^j + \beta)}{\prod_j \Gamma(n_{k-i}^j + \beta)} \frac{\Gamma(n_{k-i} + J\beta)}{\Gamma(n_{k-i} + N_i + J\beta)}
 \end{aligned}$$

GRUNDLER AND RABOSKY

243 where N_i is the total number of observations for node i and N_i^j is the subset of those
244 observations that represent utilization of resource category j . Note that for nodes with no
245 associated count data (which includes all internal nodes), this factor is equal to 1. To perform the
246 Gibbs update in practice, we calculate the conditional likelihood for each diet state according to
247 equation (7) and choose a state with probability proportional to its conditional likelihood, using
248 the sum of all conditional likelihoods as a normalizing constant. The marginal posterior
249 probability that a node is in a given diet state is simply the fraction of posterior samples where it
250 appears in that state.

251 Once a state is sampled for a node any count data associated with that node are added to
252 the set of count data generated from the sampled state. Because the Dirichlet distribution is
253 conjugate to the multinomial distribution, the posterior distribution of the multinomial
254 distribution underlying each state is also Dirichlet distributed with parameter $(n_k^1 + \beta, \dots, n_k^J +$
255 $\beta)$. During the course of updating node states we keep track of the average expected proportional
256 utilization of each resource by each state. The expected proportional utilization of resources is
257 simply the mean of the posterior distribution which is $\left(\frac{n_k^1 + \beta}{\sum_j n_k^j + J\beta}, \dots, \frac{n_k^J + \beta}{\sum_j n_k^j + J\beta}\right)$.

258 *Updating β .*—The symmetric hyperparameter β controls the shape of the Dirichlet prior
259 distribution on the latent multinomial distributions underlying each resource state. When $\beta = 1$
260 the distribution is uniform over the J -dimensional simplex of resources. When $\beta < 1$ the
261 distribution concentrates toward the corners of the simplex, and when $\beta > 1$ the distribution
262 concentrates toward the center. Because empirical datasets are typically sparse with many zeros,
263 we assume that β is uniformly distributed on the interval $(0, 1)$ and update its value using a
264 sliding window proposal mechanism. The prior and proposal ratios are 1.

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

265 *Updating α .*—When equation (6) is used to compute the likelihood the hyperparameter α
266 controls phylogenetic signal. Although it can take on any positive value the likelihood surface
267 plateaus relatively quickly as its magnitude increases and phylogenetic signal decays. By solving
268 the logarithm of (6) for the maximum likelihood estimate of α we find that

$$269 \quad \hat{\alpha} = -\frac{\log\left(\frac{fK-1}{K-1}\right)}{\log(K/(K-1)+1)}$$

270 where f is the fraction of nodes that have the same state as their ancestor. Values of $f \leq \frac{1}{K}$
271 are consistent with infinite values of α . We therefore bound α above by the value

$$272 \quad -\log\left(\frac{\lfloor N/K \rfloor + 1}{N} K^{-1}\right) / \log(K/(K-1)+1),$$
 where N is the number of nodes (not including the

273 root) in the phylogeny. We assume that α is uniformly distributed between zero and this upper
274 value and update its value using a sliding window proposal mechanism. The prior and proposal
275 ratios are 1.

276 *Updating Λ .*—When equation (4) is used to compute the likelihood the parameter Λ
277 controls phylogenetic signal. As for α , we want a reasonable upper bound for Λ because the
278 likelihood plateaus relatively quickly as the rate increases. By replacing branch lengths in
279 equation (4) with the average branch length and solving for the maximum likelihood estimate of
280 Λ under the assumption of no rate variation we find that

$$281 \quad \hat{\Lambda} = -\frac{\log\left(\frac{fK-1}{K-1}\right)}{K\bar{t}}$$

282 where \bar{t} is the average branch length. As before, values of $f \leq \frac{1}{K}$ are consistent with

283 infinite values of Λ . We therefore bound Λ above by the value $-\log\left(\frac{\lfloor N/K \rfloor + 1}{N} K^{-1}\right) / K\bar{t}$. We

GRUNDLER AND RABOSKY

284 assume that Λ is uniformly distributed between zero and this upper value and update its value
285 using a sliding window proposal mechanism. The prior and proposal ratios are 1. We use the
286 same prior distributions and proposal mechanisms detailed by Drummond and Suchard (2010)
287 for updating the number of local random clocks and the rate multipliers associated with local
288 clocks.

289

290 *Implementation*

291 Functions for fitting the model to data are provided as an R package available from
292 github.com/blueraleigh/phyr. The package includes two R functions that call compiled C
293 programs implementing the random local clocks and ultra-common mechanism models.

294

295 *Simulation study*

296 To illustrate application of the method we designed a simulation study using an empirical
297 dataset on pseudoboine snake diets (Alencar et al. 2013). Our rationale for basing simulations on
298 an empirical dataset is to ensure that properties of the data used to evaluate performance of the
299 method are consistent with real studies, especially the distribution of observations per taxon and
300 the distribution of resource specialization. Pseudoboine snakes are common members of the
301 squamate communities found in lowland rainforests of South America. Predominantly terrestrial
302 or semi-arboreal, these snakes mainly eat small mammals, lizards, and other snakes. The dataset
303 includes 606 observations of prey items from 8 prey categories for 32 species of pseudoboine
304 snakes. Per species sample sizes range from 1 to 56 observations (or 0.125-fold to seven-fold
305 coverage, where coverage is the number of observations divided by the number of resource
306 categories). We reanalyzed these data using a 33-species pseudoboine phylogeny extracted from

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

307 the posterior distribution of trees in Tonini et al. (2016). The dataset is illustrated in Figure 2.
308 The original publication coded each species with at least 8 diet observations into a set of 5
309 specialist diet states and 1 generalist diet state. Species were considered specialists if the
310 resource represented at least 70 percent of recorded prey items (as in our Figure 1). When we
311 applied the resampling procedure illustrated in Figure 1 to this empirical dataset under the
312 assumption that the original state assignments represented the “truth”, we found that in
313 approximately 20 percent of resampled datasets at least one original state (not always the same
314 state) was not present and that in about 84 percent of cases at least one species was coded
315 incorrectly (although overall coding accuracy was high, ranging from 0.77 to 1). Thus, this
316 dataset illustrates some of the concerns raised in our introduction but is also well-sampled
317 enough and shows enough variation to facilitate the estimation of separate multinomial
318 distributions.

319 Simulated datasets were generated from $K = 2, 3, 4,$ and 5 diet states using the empirical
320 sample size distribution with the original 8 food resource categories. For each K we first
321 performed Bayesian inference under the ultra-common mechanism model to estimate the
322 unobserved multinomial distributions. The estimated multinomial distributions were
323 subsequently used to simulate diet observations. For each K we simulated 20 datasets at each of
324 7 different levels of phylogenetic signal (0.1, 0.3, 0.5, 0.6, 0.7, 0.8, and 0.9) using the transition
325 probabilities in both equations (3) and (5), resulting in 560 datasets for each model and 1,120
326 datasets altogether. We defined phylogenetic signal as $p_{ii} - p_{ji}$, which ranges from 0 to 1 and
327 quantifies how much information a descendant’s state provides about the state of its ancestor
328 (Royer-Carenzi et al. 2013). Using equation (5) for transition probabilities results in phylogenetic
329 signal equal to $\left(\frac{1}{K/(K-1)+1}\right)^\alpha$. We used this result to calculate the value of α for each simulation.

GRUNDLER AND RABOSKY

330 When equation (3) is used for transition probabilities each branch has a unique phylogenetic
331 signal. Because phylogenetic signal is a convex function of branch length, the average
332 phylogenetic signal of all branches is greater than or equal to the phylogenetic signal of the
333 average branch, which is $e^{-K\Lambda\bar{t}}$. We used the phylogenetic signal of the average branch to
334 calculate the value of Λ for each simulation, which we applied to all branches (i.e., datasets did
335 not include random local clock variation). Interestingly, for a given branch length (measured as
336 expected number of state changes) phylogenetic signal with equation (5) is always greater than
337 phylogenetic signal with equation (3), suggesting that estimating the rate of evolution trades off
338 with estimating ancestral node states (Gascuel and Steel 2018). For each simulated dataset we
339 ran a set of Markov chains with 1, 2, ..., up to $K+3$ diet states. Each chain was run for 160,000
340 iterations after a burnin of 30,000 iterations, sampling every 128 iterations to yield 1,250
341 posterior samples.

342

343 *Determining the number of resource states*

344 Because the model does not include a process for generating the number of states, we
345 must perform analyses across multiple values of K and apply an *a posteriori* inference procedure
346 to choose between them. A similar problem is encountered when trying to infer the number of
347 demes from multi-locus genotype data with the program STRUCTURE (Pritchard et al. 2000).
348 Our approach is to choose the largest value of K for which all states are unambiguously assigned
349 to at least one terminal taxon. Specifically, by examining the terminal nodes of the phylogeny we
350 can determine the maximum marginal posterior probability assigned to each state across all
351 terminal nodes. Call the smallest of these q_K . A low value of q_K implies that at least one state is
352 not assigned to any individuals (terminal taxa) with high probability. As K varies from low to

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

353 high there will come a point when additional states become redundant to previous states. When
354 that occurs, q_K will drop well below 1. In other words, a point will be reached when at least one
355 state is never unambiguously assigned to at least one terminal node. To choose the best value for
356 K we first identify the K with the steepest drop in q_K with respect to q_{K-1} . $K^* = K - 1$ then
357 becomes a candidate best choice. If $q_{K^*} = 1$ (or nearly so) we treat K^* as the best choice.
358 Otherwise, we keep setting $K^* = K^* - 1$ until $q_{K^*} = 1$. This procedure, which we call the q_K
359 rule, is illustrated in Figure 3 on the empirical dataset.

360

361 *Assessing model adequacy*

362 In practice, the number of states identified by the q_K rule will depend on the underlying
363 data and may change as data are updated. This is because the model may be unable to distinguish
364 truly different samples as having arisen from separate distributions if sample sizes are small.
365 Therefore, empirical applications of the method require a way of assessing how well the inferred
366 multinomial distributions explain the empirical resource observations. Our approach is to
367 compute a per-species adequacy score that measures the similarity of each taxon, with respect to
368 sampled observations, to other taxa assigned to the same state. The procedure we describe
369 effectively measures the compositional heterogeneity of the model-inferred states with respect to
370 sampled diet observations. If the model is fully adequate and accurately describes patterns of
371 resource use, then all species assigned to a given state will have similar sampled diets (e.g., the
372 observed data). However, some species may be assigned to states even where they have
373 dissimilar sampled observations from other species in the same state, reflecting overdispersed
374 diet distributions for the state. Such overdispersion might arise if species are assigned to the
375 “wrong” state. These incorrect state assignments might be preferred under the model if a species

GRUNDLER AND RABOSKY

376 does not have enough observations to provide information about the existence of a distinct state,
377 and so the species is conservatively assigned to the immediately ancestral diet state. To compute
378 the adequacy score, we first draw a state for each terminal node from its array of marginal
379 posterior probabilities computed using equation (7) (which utilizes count data and evolutionary
380 history). Then we visit each terminal node in turn and perform the following exercise. Given the
381 configuration of states for all terminal nodes, we identify the set of count data generated by the
382 state of the current node and compute the likelihood of these data using equation (1) (which only
383 utilizes count data). Next, we compute the likelihood of these same data, using equation (1)
384 together with equation (2), but assuming that they were generated from two states by placing the
385 current node in its own unique state. Finally, we take the negative log likelihood ratio of these
386 two values. We repeat this procedure for a thousand independent configurations and record the
387 average negative log likelihood ratio (adequacy score) for each terminal node. Large negative
388 adequacy scores highlight terminal nodes that were assigned to a state whose other members
389 have a strongly dissimilar pattern of resource utilization. This procedure is illustrated in Figure 4
390 on the empirical dataset (see also Figure S4).

391

392 RESULTS

393 Overall, the q_K rule correctly identified the number of resource states in 492 of 560
394 simulations from the ultra-common mechanism model (Fig. 5). In the 68 cases where the method
395 incorrectly identified the number of states, it underestimated the number of states by one (61
396 instances), two (4 instances), and three states (2 instance) and overestimated the number of states
397 by one state in one instance. When the q_K rule was used with the random local clocks model it
398 correctly identified the number of states in 475 of 560 simulations (Fig. S1). In the 85 cases

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

399 where the method incorrectly identified the number of states, it underestimated the number of
400 states by one (77 instances) and two states (8 instances). Failure to correctly identify the number
401 of states commonly occurs when the number of observations generated by a state is small
402 relative to the number of observations from other states. This happens when the terminal nodes
403 representing a state have poorly sampled diets causing the state to be subsumed into the state of
404 near relatives.

405 Estimation of latent multinomial probabilities underlying resource states was highly
406 accurate (Figs. 5, S1, S5). In simulations from both the ultra-common mechanism model and the
407 random local clocks model, the correlation between true and estimated probabilities ranged from
408 0.92 to 0.999 with a mean of 0.99. Such high values occur because even when the number of
409 states is underestimated, the number of sampled observations from the missing states are few
410 enough in number that they do not appreciably alter the estimated proportions of the other states
411 (Fig S6).

412 Across all levels of phylogenetic signal, and with simulations from both the ultra-
413 common mechanism model and the random local clocks model, the method consistently
414 classified greater than ninety-percent of terminal nodes in the correct state. For internal nodes,
415 reconstruction accuracy was comparable to reconstruction accuracy for terminal nodes at the
416 highest level of phylogenetic signal but, in accordance with expectation, decayed toward the
417 expected accuracy of a random guess as phylogenetic signal declined toward zero (Figs. 6, S2).

418 This behavior was mirrored in the posterior distributions of α and Λ (Figs. 7, S3). When
419 phylogenetic signal was low, posterior distributions of rate estimates were diffuse and resembled
420 their uniform prior distributions except that they were shifted away from the lower bound. As

GRUNDLER AND RABOSKY

421 phylogenetic signal increased, posterior distributions of rate estimates concentrated around the
422 true values used to generate simulated datasets.

423

424 DISCUSSION

425 We developed a comparative method for macroevolutionary analysis of multivariate
426 count data. The method is general and may be applied to any data expressible as a set of
427 observational counts from different categories. Such datatypes arise frequently in community
428 ecology and behavior. Potential applications include the comparative analysis of diet, foraging
429 behavior, activity patterns, and habitat preferences. The method is similar to standard
430 continuous-time Markov chain models of phenotypic evolution but differs in several important
431 respects. First, the number of states in the model and the states themselves are unobserved and
432 must be estimated from empirical data on resource use. Second, each state is an unobserved
433 multinomial distribution rather than a categorical variable. This latter property enables
434 researchers to model ecological traits that show intraspecific variation and to account for
435 uncertainty in the state assignments of terminal taxa that arises from the effects of sampling
436 variation.

437 Simulations revealed that the new method is generally able to determine the correct
438 number of states and that it provides accurate estimates of the underlying (unobserved)
439 multinomial distributions, both for terminal taxa as well as internal nodes. We designed
440 simulations around empirical patterns of resource use in a dataset on snake diets (Alencar et al.
441 2013). Therefore, caution is warranted in generalizing the good performance observed in the
442 current study to other datasets. In particular, performance of the model will depend on the
443 idiosyncrasies of individual datasets, including the distribution of sample sizes and the

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

444 distribution of overlaps in resource use among species. We expect that states represented by few
445 observations will be difficult to infer, especially if those states show appreciable overlap with
446 other states.

447 Our empirical analysis identified at least two feeding modalities among the set of species
448 Alencar et al. (2013) recognize as “generalists”: species that feed predominantly on snakes but
449 that regularly eat lizards and mammals, and species that feed predominantly on mammals and
450 lizards. Ancestral state estimates strongly suggest that each of these feeding modalities arose
451 from a more specialized diet comprised almost entirely of lizards. This is in contrast to the
452 results of Alencar et al. (2013), which imply that nearly all origins of specialized feeding
453 modalities occurred from a generalist ancestor, although direct comparison of results is made
454 difficult by the use of different phylogenies in the two studies.

455 As currently implemented, the approach described here does not directly model gains and
456 losses or substitutions of different resources. Indeed, no resource is ever truly absent from the
457 reconstructed states (although its proportional representation may approach zero as β becomes
458 small). This contrasts with biogeographic-type models that explicitly model resource use
459 expansions, contractions, and substitutions (e.g. see Hardy (2017) for application of Ree and
460 Smith’s (2008) dispersal-extinction-cladogenesis model to binary encoded diet data). Although
461 these types of changes are implicit in the sequence of reconstructed states derived from the
462 model, future studies might want to explore how to combine more complex evolutionary models
463 with the current model for count data. Nonetheless, the advantage of a simple evolutionary
464 model is that it has broad scope. It would be possible, for example, to apply our method to
465 continuous characters by keeping the same evolutionary model but changing the model for
466 observations from a Dirichlet-multinomial to a multivariate-normal distribution, which could

GRUNDLER AND RABOSKY

467 then be applied to other data types used for quantifying resource use such as stable isotope ratios
468 of carbon and nitrogen.

469 One challenge for comparative methods is their limited ability to model ecological
470 phenotypes that cannot be neatly summarized by a single value (Hardy & Linder 2005; Hardy
471 2006). Recent years have seen progress in this direction for continuous traits, including models
472 that accommodate intraspecific variation, function-valued traits, and other non-gaussian data
473 (Ives et al. 2007; Felsenstein 2008; Evans et al. 2009; Jones & Moriarty 2013; Goolsby 2015;
474 Quintero et al. 2015). The general approach developed here, where each state is a multinomial
475 distribution rather than a categorical variable, extends this progress to traits like diet and habitat
476 that are typically treated as categorical variables. By placing an emphasis on individual natural
477 history observations, the method draws attention to the central role such observations play in
478 evolutionary biology (Greene 1986) and to the many remaining opportunities for developing
479 comprehensive ecological databases that advance our understanding of biodiversity (Hortal et al.
480 2015).

481

482 SUMMARY

483 We described a novel methodological framework for studying the evolutionary dynamics
484 of complex ecological traits on phylogenetic trees. Previous approaches to this problem have
485 assumed that ecological states are categorical variables and that species are monomorphic for
486 particular states. We relaxed this assumption through the use of a hidden Markov model that
487 treats ecological states as unobserved probability distributions from which observed data are
488 sampled. Although our method is designed for the analysis of multivariate count data, we suggest

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

489 that the approach of treating states as hidden probability distributions has wide applicability and
490 will greatly facilitate the comparative analysis of novel sources of ecological data.

491

492 ACKNOWLEDGEMENTS

493 This research was supported in part through computational resources and services provided by
494 Advanced Research Computing at the University of Michigan, Ann Arbor. This material is based
495 upon work supported by the National Science Foundation Graduate Research Fellowship and by
496 the University of Michigan Department of Ecology and Evolutionary Biology Block Grant.

497

498 LITERATURE CITED

- 499 Alencar L.R.V., Gaiarsa M.P., Martins M. 2013. The evolution of diet and microhabitat use in
500 pseudoboine snakes. *South American Journal of Herpetology*. 8:60–66.
- 501 Bassar R.D., Marshall M.C., López-Sepulcre A., Zandona E., Auer S.K., Travis J., Pringle C.M.,
502 Flecker A.S., Thomas S.A., Fraser D.F., Reznick D.N. 2010. Local adaptation in
503 Trinidadian guppies alters ecosystem processes. *PNAS*. 107:3616–3621.
- 504 Beaulieu J.M., O’Meara B.C. 2016. Detecting hidden diversification shifts in models of trait-
505 dependent speciation and extinction. *Systematic Biology*. 65:583–601.
- 506 Beaulieu J.M., O’Meara B.C., Donoghue M.J. 2013. Identifying hidden rate changes in the
507 evolution of a binary morphological character: the evolution of plant habit in Campanulid
508 angiosperms. *Systematic Biology*. 62:725–737.
- 509 Blei D.M., Ng A.Y., Jordan M.I. 2003. Latent Dirichlet Allocation. *Journal of Machine Learning*
510 *Research*. 3:993–1022.
- 511 Burin G., Kissling W.D., Guimarães P.R., Sekercioglu C.H., Quental T.B. Omnivory in birds is a
512 macroevolutionary sink. *Nature Communications*. 7:11250.
- 513 Caetano D.S., Beaulieu J.M., O’Meara B.C. 2018. Hidden state models improve state-dependent
514 diversification approaches, including biogeographical models. *Evolution*. 72:2308–2324.
- 515 Cantalapiedra J.L., FitzJohn R.G., Kuhn T.S., Hernández Fernández M., DeMiguel D., Azanza
516 B., Morales J., Mooers A.Ø. 2014. Dietary innovations spurred the diversification of
517 ruminants during the Cenozoic. *Proceedings of the Royal Society B: Biological*
518 *Sciences*. 281:20132746.

GRUNDLER AND RABOSKY

- 519 Colston T.J., Costa G.C., Vitt L.J. 2010. Snake diets and the deep history hypothesis. *Biological*
520 *Journal of the Linnean Society*. 101:476–486.
- 521 Davis A.M., Unmack P.J., Vari R.P., Betancur-R R. 2016. Herbivory promotes dental
522 disparification and macroevolutionary dynamics in grunters (Teleostei: Terapontidae), a
523 freshwater adaptive radiation. *The American Naturalist*. 187:320–333.
- 524 Drummond A.J., Suchard M.A. 2010. Bayesian random local clocks, or one rate to rule them all.
525 *BMC Biology*. 8:114.
- 526 Ehrlich P.R., Raven P.H. 1964. Butterflies and plants: a study in coevolution. *Evolution*. 18:586–
527 608.
- 528 Evans M.E.K., Smith S.A., Flynn R.S., Donoghue M.J. 2009. Climate, niche evolution, and
529 diversification of the “bird-cage” evening primroses (*Oenothera*, Sections *Anogra* and
530 *Kleinia*). *The American Naturalist*. 173:225–240.
- 531 Felsenstein J. 2008. Comparative methods with sampling error and within-species variation:
532 contrasts revisited and revised. *The American Naturalist*. 171:713–725.
- 533 Felsenstein J. 2012. A comparative method for both discrete and continuous characters using the
534 threshold model. *The American Naturalist*. 179:145–156.
- 535 Forister M.L., Dyer L.A., Singer M.S., Stireman III J.O., Lill J.T. 2012. Revisiting the evolution
536 of ecological specialization, with emphasis on insect-plant interactions. *Ecology*. 93:981–
537 991.
- 538 Futuyma D.J., Moreno G. 1988. The evolution of ecological specialization. *Annual Review of*
539 *Ecology and Systematics*. 19:207–233.
- 540 Gajdzik L., Aguilar-Medrano R., Frédéricich B. 2019. Diversification and functional evolution of
541 reef fish feeding guilds. *Ecology Letters*. 22:572–582.
- 542 Gascuel O., Steel M. 2018. A Darwinian uncertainty principle. bioRxiv.
543 <https://doi.org/10.1101/506535>
- 544 Gillespie R. 2004. Community assembly through adaptive radiation in Hawaiian spiders.
545 *Science*. 303:356–359.
- 546 Givnish T.J., Barfuss M.H.J., Van Ee B., Riina R., Schulte K., Horres R., Gonsiska P.A., Jabaily
547 R.S., Crayn D.M., Smith J.A., Winter K., Brown G.K., Evans T.M., Holst B.K., Luther
548 H., Till W., Zizka G., Berry P.E., Sytsma K.J. 2014. Adaptive radiation, correlated and
549 contingent evolution, and net diversification in Bromeliaceae. *Molecular Phylogenetics*
550 *and Evolution*. 71:55–78.
- 551 Goolsby E.W. 2015. Phylogenetic Comparative Methods for Evaluating the Evolutionary
552 History of Function-Valued Traits. *Systematic Biology*. 64:568–578.

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

- 553 Greene H.W. 1986. Natural history and evolutionary biology. In: Feder M.E., Lauder G.V.,
554 editors. Predator–prey relationships: perspectives and approaches from the study of lower
555 vertebrates. University of Chicago Press. p. 198.
- 556 Hardy C.R. 2006. Reconstructing ancestral ecologies: challenges and possible solutions.
557 Diversity and Distributions. 12:7–19.
- 558 Hardy C.R., Linder H.P. 2005. Intraspecific variability and timing in ancestral ecology
559 reconstruction: a test case from the cape flora. Systematic Biology. 54:299–316.
- 560 Hardy N.B. 2017. Do plant-eating insect lineages pass through phases of host-use generalism
561 during speciation and host switching? Phylogenetic evidence. Evolution. 71:2100–2109.
- 562 Hardy N.B., Otto S.P. 2014. Specialization and generalization in the diversification of
563 phytophagous insects: tests of the musical chairs and oscillation hypotheses. Proceedings
564 of the Royal Society B: Biological Sciences. 281:20132960.
- 565 Harmon, L. J., Matthews B., Des Roches S., Chase J.M., Shurin J.B., Schluter D. 2009.
566 Evolutionary diversification in stickleback affects ecosystem functioning. Nature.
567 458:1167–1170.
- 568 Hastings W.K. 1970. Monte Carlo Sampling Methods Using Markov Chains and Their
569 Applications. Biometrika. 57:97–109.
- 570 Hortal J., de Bello F., Diniz-Filho J.A.F., Lewinsohn T.M., Lobo J.M., Ladle R.J. 2015. Seven
571 shortfalls that beset large-scale knowledge of biodiversity. Annual Review of Ecology,
572 Evolution, and Systematics. 46:523–549.
- 573 Huelsenbeck J.P., Ané C., Larget B., Ronquist F. 2008. A Bayesian perspective on a non-
574 parsimonious parsimony model. Systematic Biology. 57:406–419.
- 575 Ives A.R., Midford P.E., Garland Jr T. 2007. Within-species variation and measurement error in
576 phylogenetic comparative methods. Systematic Biology. 56:252–270.
- 577 Janz N., Nyblom K., Nylin S. 2001. Evolutionary dynamics of host-plant specialization: a case
578 study of the tribe Nymphalini. Evolution. 55:783–796.
- 579 Jones N.S., Moriarty J. 2013. Evolutionary inference or function-valued traits: Gaussian process
580 regression on phylogenies. Journal of the Royal Society Interface. 10:20120616.
- 581 Kelley S.T., Farrell B.D. 1998. Is specialization a dead end? The phylogeny of host use in
582 *Dendroctonus* bark beetles (Scolytidae). Evolution. 52:1731–1743.
- 583 Losos J.B., Leal M., Glor R.E., de Queiroz K., Hertz P.E., Rodríguez Schettino L., Chamizo Lara
584 A., Jackman T.R., Larson A. 2003. Niche lability in the evolution of a Caribbean lizard
585 community. Nature. 424:542–545.

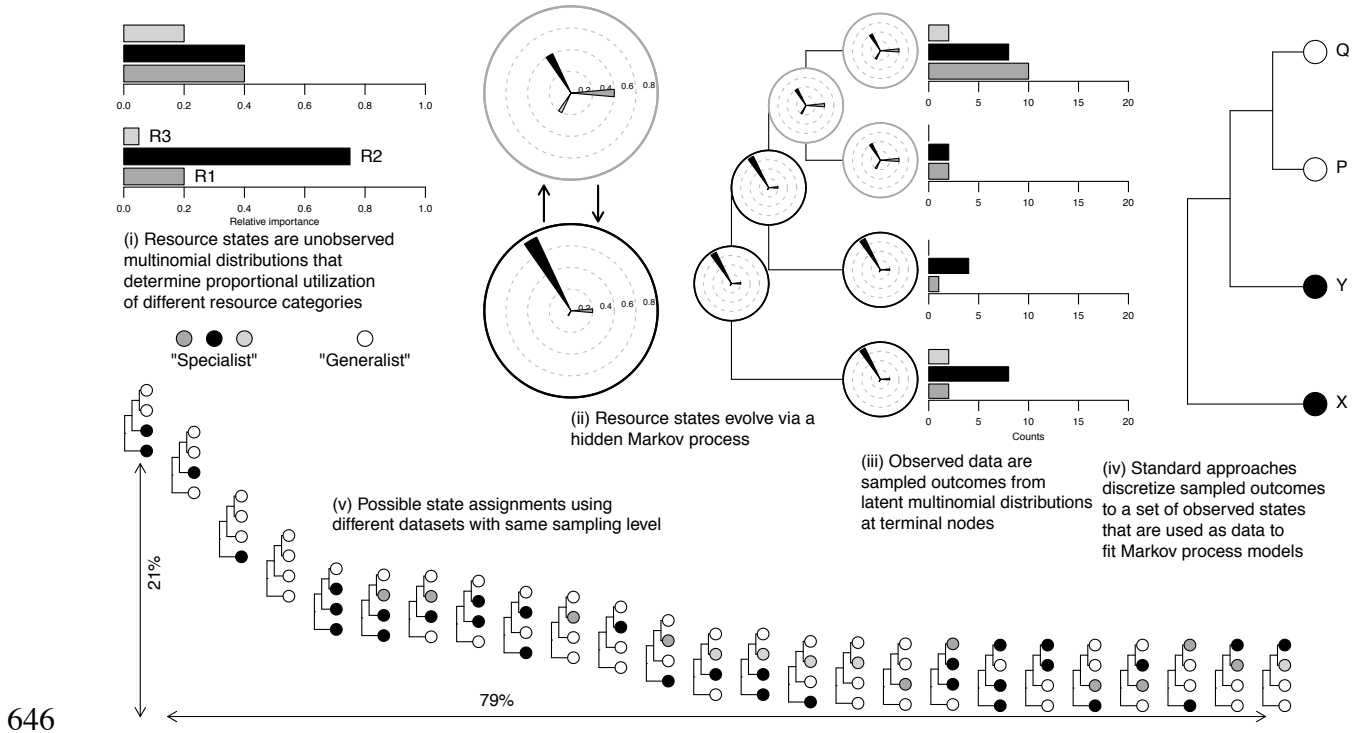
GRUNDLER AND RABOSKY

- 586 MacArthur R.H. 1958. Population ecology of some warblers of northeastern coniferous forests.
587 Ecology. 39:599–619.
- 588 Marazzi B., Ané C., Simon M.F., Delgado-Salinas A., Luckow M., Sanderson M.J. 2012.
589 Locating evolutionary precursors on a phylogenetic tree. Evolution. 66:3918–3930.
- 590 Martin C.H., Wainwright P.C. 2011. Trophic novelty is linked to exceptional rates of
591 morphological diversification in two adaptive radiations of *Cyprinodon* pupfish.
592 Evolution. 65:2197–2212.
- 593 Miller E.T., Wagner S.K., Harmon L.J., Ricklefs R.E. 2017. Radiating despite a lack of
594 character: ecological divergence among closely related, morphologically similar
595 honeyeaters (Aves: Meliphagidae) co-occurring in arid Australian environments. The
596 American Naturalist. 189:E14–E30.
- 597 Mitter C., Farrell B., Wiegmann B. 1988. The phylogenetic study of adaptive zones: has
598 phytophagy promoted insect diversification? The American Naturalist. 132:107–128.
- 599 Nosil P. 2002. Transition rates between specialization and generalization in phytophagous
600 insects. Evolution. 56:1701–1706.
- 601 O’Meara B.C. 2012. Evolutionary inferences from phylogenies: a review of methods. Annual
602 Review of Ecology and Systematics. 43:267–285.
- 603 Price S.A., Hopkins S.S.B., Smith K. K., Roth V.L. 2012. Tempo of trophic evolution and its
604 impact on mammalian diversification. PNAS. 109:7008–7012.
- 605 Pritchard J.K., Stephens M., Donnelly P. 2000. Inference of Population Structure Using
606 Multilocus Genotype Data. Genetics. 155:945–959.
- 607 Quintero I., Keil P.K., Jetz W., Crawford F.W. 2015. Historical biogeography using species
608 geographical ranges. Systematic Biology. 64:1059–1073.
- 609 Ree R.H., Smith S.A. 2008. Maximum likelihood inference of geographic range evolution by
610 dispersal, local extinction, and cladogenesis. Systematic Biology. 57:4–14.
- 611 Revell L.J. 2014. Ancestral character estimation under the threshold model from quantitative
612 genetics. Evolution. 68:743–759.
- 613 Robinson G.S. 1999. HOSTS - a database of the hostplants of the world’s Lepidoptera. Nota
614 Lepidopterologica. 22:35–47.
- 615 Royer-Carenzi M., Pontarotti P., Didier G. 2013. Choosing the best ancestral character state
616 reconstruction method. Mathematical Biosciences. 242:95–109.
- 617 Shine R. 1994. Allometric patterns in the ecology of Australian snakes. Copeia. 1994:851–867.

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

- 618 Steel M. 2011. Can we avoid “SIN” in the house of “no common mechanism”? *Systematic*
619 *Biology*. 60:96–109.
- 620 Tonini J., Beard K.H., Barbosa Ferreira R., Jetz W., Pyron R.A. 2016. Fully-sampled
621 phylogenies of squamates reveal evolutionary patterns in threat status. *Biological*
622 *Conservation*. 204:23–31.
- 623 Tuffley C., Steel M. 1997. Links between maximum likelihood and maximum parsimony under a
624 simple model of site substitution. *Bulletin of Mathematical Biology*. 59:581–607.
- 625 Vitt L.J., Vangilder L.D. 1983. Ecology of a snake community in northeastern Brazil. *Amphibia-*
626 *Reptilia*. 4:273–296.
- 627 Vrba E.S. 1987. Ecology in relation to speciation rates: some case histories of Miocene-Recent
628 mammal clades. *Evolutionary Ecology*. 1:283–300.
- 629 Yin J., Wang J. 2014. A dirichlet multinomial mixture model-based approach for short text
630 clustering. *Proceedings of the 20th ACM SIGKDD international conference on*
631 *Knowledge discovery and data mining - KDD '14*.:233–242.
- 632
- 633
- 634
- 635
- 636
- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645

GRUNDLER AND RABOSKY



646

647 **Figure 1.** Distribution and representation of multivariate ecological phenotypes (i, ii, iii), data as
 648 sampled by researchers (iii), and sampled states as typically represented by univariate categorical
 649 traits (iv, v). Loss and distortion of information associated with complex phenotypes motivates
 650 the development of the Dirichlet-multinomial model described in this article. (i) True resource
 651 states are unobserved multinomial distributions that determine the proportional utilization of
 652 three dietary resource categories by four species. (ii) The resource state of a species is the
 653 outcome of evolution via a hidden Markov process where the hidden states are the unobserved
 654 multinomial distributions. Here, the multinomial distributions from (i) are represented as rose
 655 plots: the direction of a spoke identifies the resource category and the length of a spoke is equal
 656 to the proportional utilization of that category. The phylogeny depicts the true evolutionary
 657 history of change. (iii) Observed data are sampled outcomes from these latent multinomial
 658 distributions. (iv) Sampled outcomes are projected onto a set of resource states. Here, a species is
 659 considered a “specialist” on a particular category if the sampled proportion of the category

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

660 exceeds 0.7. Otherwise, it is considered a “generalist”. In this case, the dataset and cutoff value
661 align to match each species with its correct modal resource category. (v) The same state
662 assignment process is repeated with different datasets while keeping the sample sizes for each
663 species identical to (iii). State assignments are sorted along the x -axis according to their
664 frequency of occurrence in 1,000 independent datasets. Datasets were generated by sampling
665 from the two multinomial distributions in (i). Note that the procedure correctly matches all
666 species with their modal food resource in a minority of cases and results in a variable number of
667 states across datasets. The implication for macroevolutionary studies using this state assignment
668 procedure is that we cannot be certain whether state assignments are reflective of true patterns of
669 resource use or are merely the expected outcome of sampling variation.

670

671

672

673

674

675

676

677

678

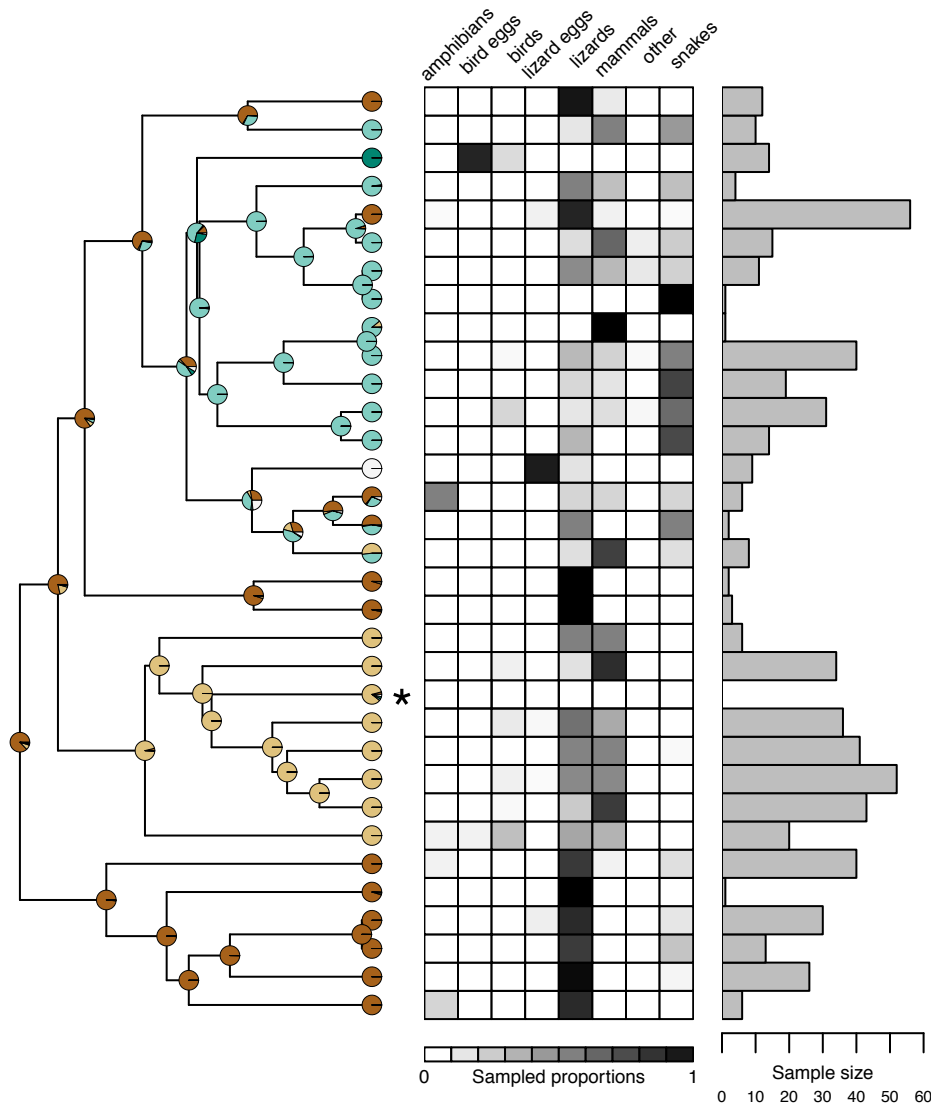
679

680

681

682

GRUNDLER AND RABOSKY



683

684 **Figure 2.** Summary attributes of the snake dietary dataset used to parameterize the simulation
685 study, including phylogeny of pseudoboine snakes (left), relative prey frequencies (middle), and
686 total numbers of food observations per snake species (right). Dark colors in the prey frequency
687 matrix indicate higher sampled proportions of a particular prey item in a given diet. The
688 phylogeny and sampled diet observations were used to infer the “true” (unobserved) diet states
689 and their evolutionary history on the phylogeny (colored circles). Each color in a pie chart is a
690 specific model-inferred state and the size of a slice represents the marginal posterior probability

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

691 that a node belongs to that state. Note that the terminal node marked with an asterisk is missing
692 data. It is treated like an internal node and information about its probable diet state is drawn only
693 from what the model has learned about the states of its neighbors and the likelihood of
694 evolutionary change. Here, the Dirichlet-multinomial model inferred 5 states, corresponding to 3
695 specialist (> 70% specificity for a single prey group) and 2 generalist diets. Note that the diet
696 states are not observed directly, even at the tips of the tree; rather, all observed data are assumed
697 to be sampled from a set of unknown multinomial distributions.

698

699

700

701

702

703

704

705

706

707

708

709

710

711

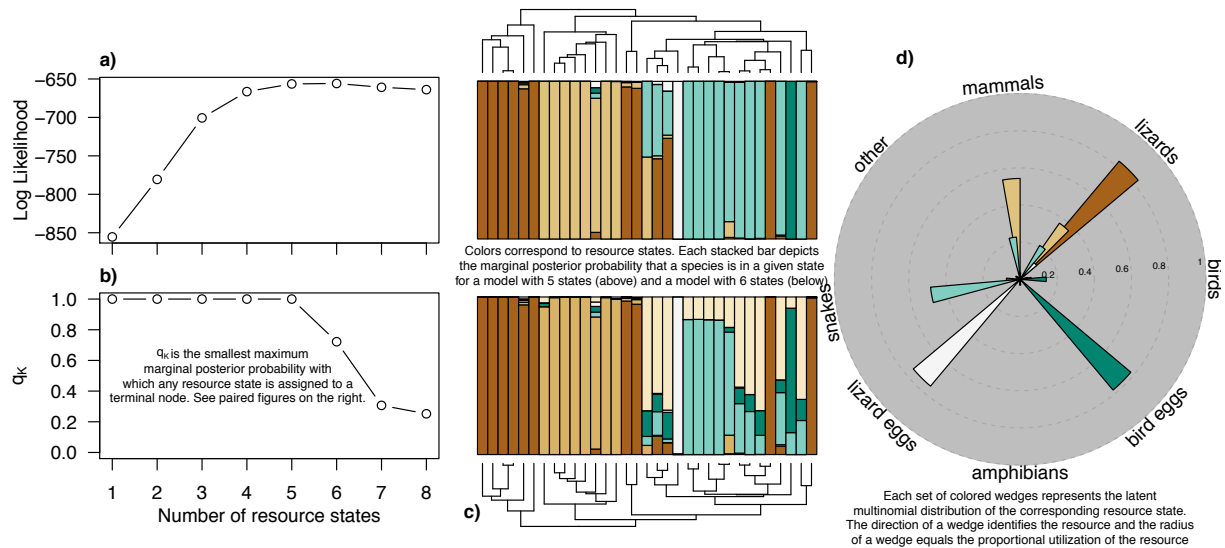
712

713

714

715

GRUNDLER AND RABOSKY



716

717 **Figure 3.** Illustration of the *a posteriori* criterion for determining the number of states in the
718 model. Panel (a) shows the average log likelihood of the empirical data as a function of the
719 number of diet states. Panel (b) depicts how q_K , the smallest maximum marginal posterior
720 probability with which a state is assigned to terminal taxa, changes as a function of the number
721 of states. Inspection of the marginal posterior probabilities reveals that the sixth state is never
722 unambiguously assigned to a terminal node (panels b and c). For this reason, a model with five
723 resource states is considered optimal. The proportional utilization of different food resources by
724 those five states is illustrated by the rose plot in panel (d).

725

726

727

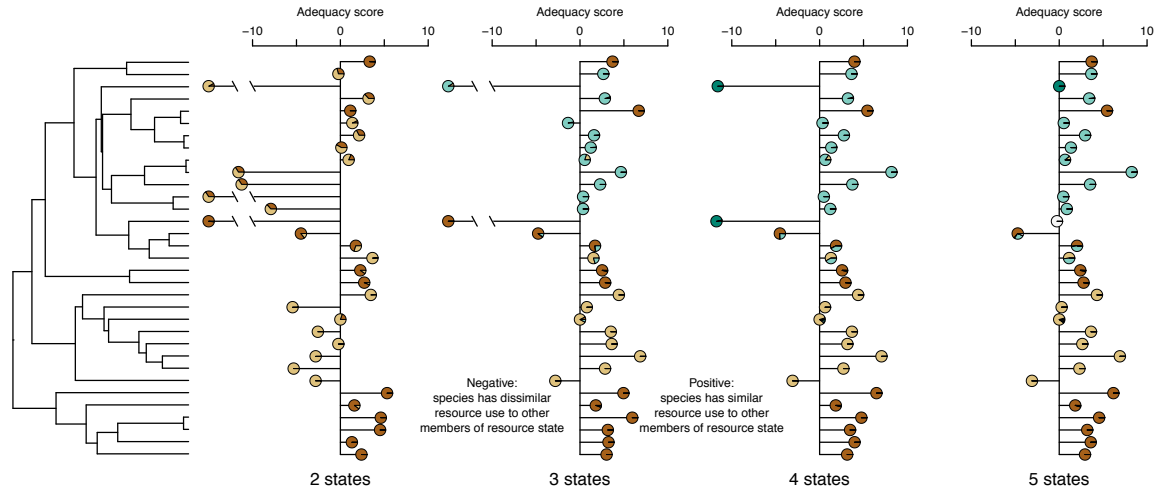
728

729

730

731

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA



732

733 **Figure 4.** Goodness of fit of inferred multinomial distributions to sampled observations in the

734 empirical snake diet dataset when the inference model assumes $K = 2, 3, 4,$ and 5 states.

735 Adequacy scores are computed on a per-species basis; negative scores imply that sampled

736 observations for a species are different from other taxa assigned to the same state. Such negative

737 scores might arise when a species has sampled diet that is distinct from that of close relatives but

738 where insufficient data (few observations) are available to inform the model as to the existence

739 of a new and distinct dietary state. Each line segment is a species in the empirical dataset. Each

740 pie chart depicts the marginal posterior probabilities that a species is in a given diet state. Thus,

741 fitting a model with only 2 states results in a third of the species sharing a diet state with

742 ecologically dissimilar species. In a model with 5 states nearly all species in the empirical dataset

743 share a similar pattern of resource utilization with the other species assigned to the same diet

744 state. The two species with negative scores have a relatively high sampled percentage of

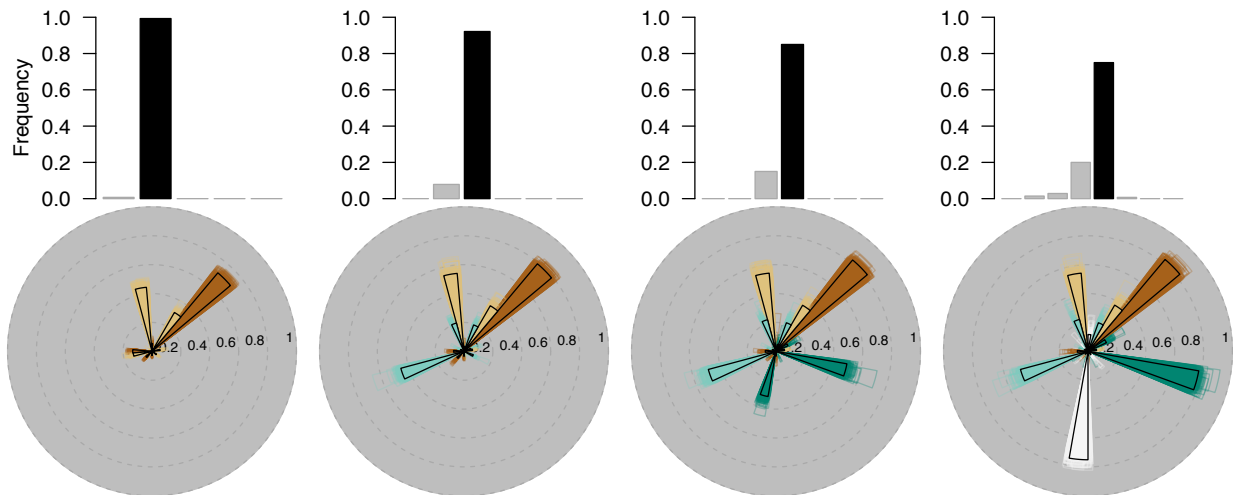
745 amphibian prey items or a relatively broad diet compared to other species in the dataset but only

746 modest sample sizes (compare to Fig. 2).

747

748

GRUNDLER AND RABOSKY



749

750 **Figure 5.** Performance of the method at identifying the true number of states (top) and the latent

751 multinomial distributions underlying states (bottom). Each column summarizes the results of a

752 set of 140 simulated datasets where the number of resource states varied from 2 (leftmost

753 column) to 5 (rightmost column). *Top:* each bar chart shows the frequency of the number of

754 resource states estimated using the q_K rule discussed in the main text; black bars are the number

755 of states in the generating model. The method correctly identifies the number of generating states

756 in most cases. *Bottom:* each set of solid colored spokes represents the latent multinomial

757 distribution underlying a resource state in the generating model. The direction of a spoke

758 identifies its corresponding resource and its length equals the proportional utilization of that

759 resource. Directions are slightly offset between diets states so that spokes with contacting edges

760 represent the same resource. Spokes with colored outlines but unfilled centers are the model

761 estimated multinomial distributions for the corresponding state. The method correctly identifies

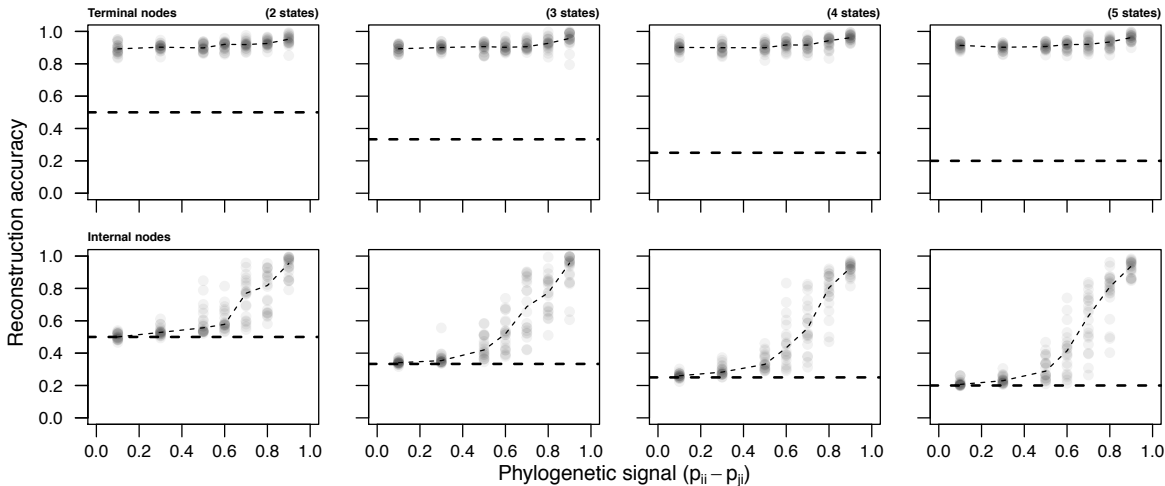
762 the major and minor resources of each state in most cases.

763

764

765

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA



766

767 **Figure 6.** Accuracy of model-inferred resource state assignments for terminal and internal nodes

768 across varying levels of phylogenetic signal. Each point is the reconstruction accuracy for a

769 single simulation. Reconstruction accuracy is the probability that a randomly selected node was

770 classified in the correct resource state. The horizontal dashed lines correspond to the expected

771 reconstruction accuracy of a random guess. The reconstruction accuracy of terminal nodes is

772 relatively constant across levels of phylogenetic signal. The reconstruction accuracy of internal

773 nodes is comparable to the accuracy of terminal nodes when phylogenetic signal is high and, as

774 expected, decays toward the random expectation as phylogenetic signal declines. This decline is

775 expected because, with low signal, the tip data provide little information about states at internal

776 nodes. Light weight dashed lines connect the median accuracy at each level of phylogenetic

777 signal.

778

779

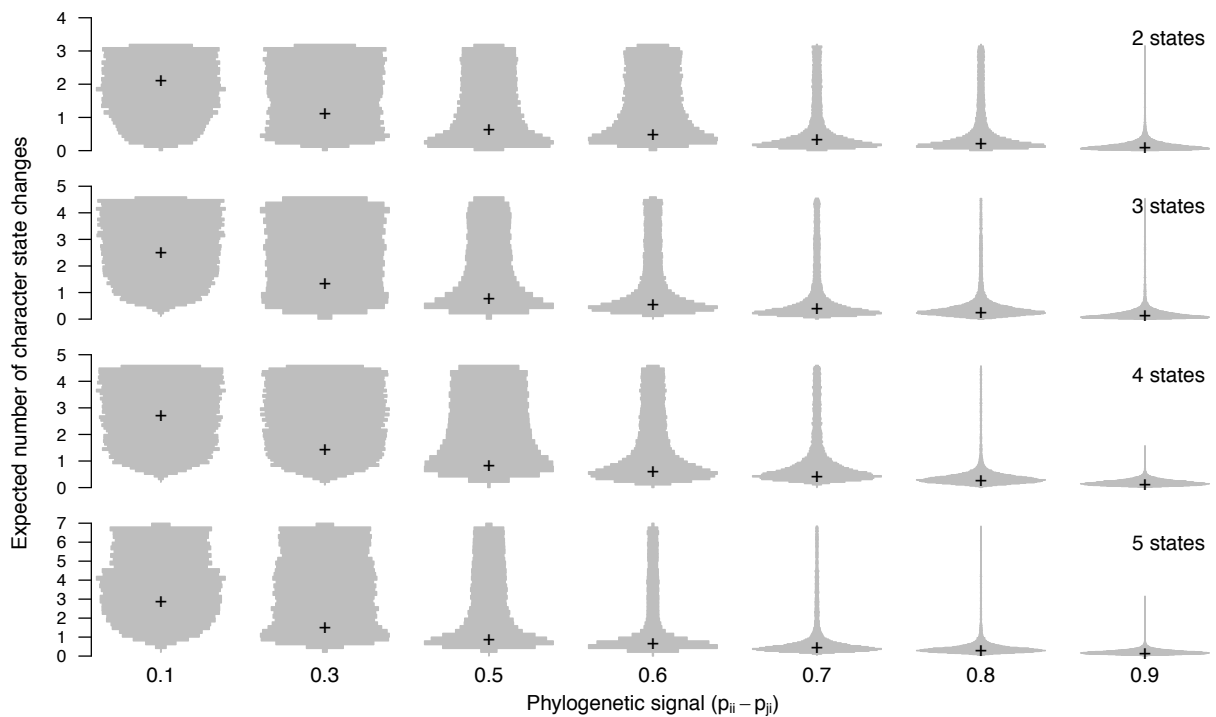
780

781

782

783

GRUNDLER AND RABOSKY



784

785 **Figure 7.** Posterior distributions of the expected number of character state changes per branch
786 for simulated datasets with $K = 2, 3, 4,$ and 5 states. Each violin plot summarizes the posterior
787 distributions of 20 simulations at each of 7 different levels of phylogenetic signal. The black
788 hatch mark within each violin is the expected number of changes per branch that was used to
789 generate simulated datasets. The width of a violin measures how frequently estimated values
790 appear in the posterior distributions. Phylogenetic signal increases as the expected number of
791 changes decreases. At low phylogenetic signals posterior distributions largely recapitulate their
792 uniform prior distributions but as phylogenetic signal increases they become concentrated around
793 the true values.

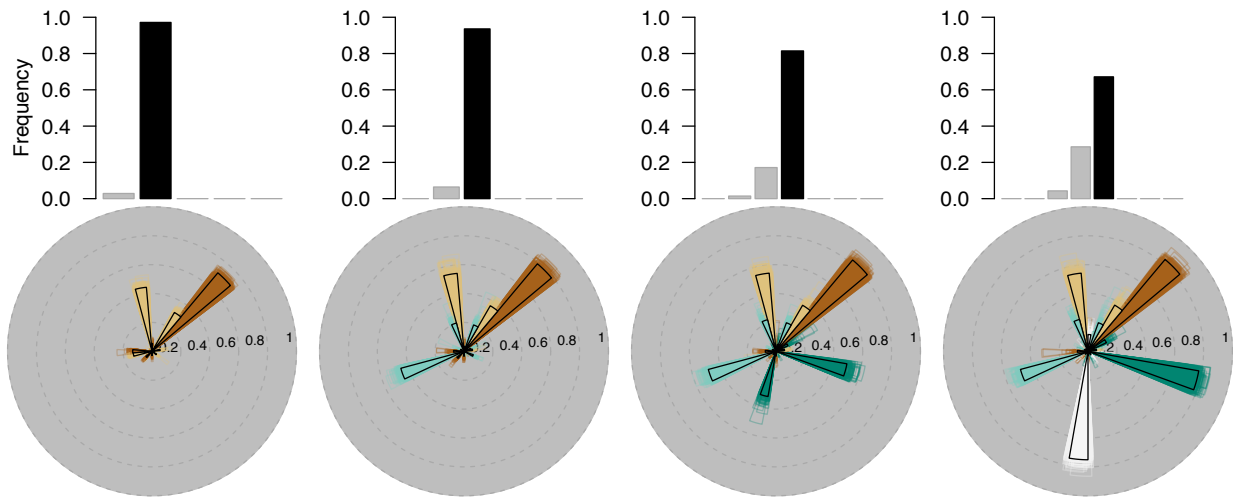
794

795

796

797

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA



798

799 **Figure S1.** As for Figure 5 in the main text except that simulations were made using transition

800 probabilities from equation (3) rather than from equation (5).

801

802

803

804

805

806

807

808

809

810

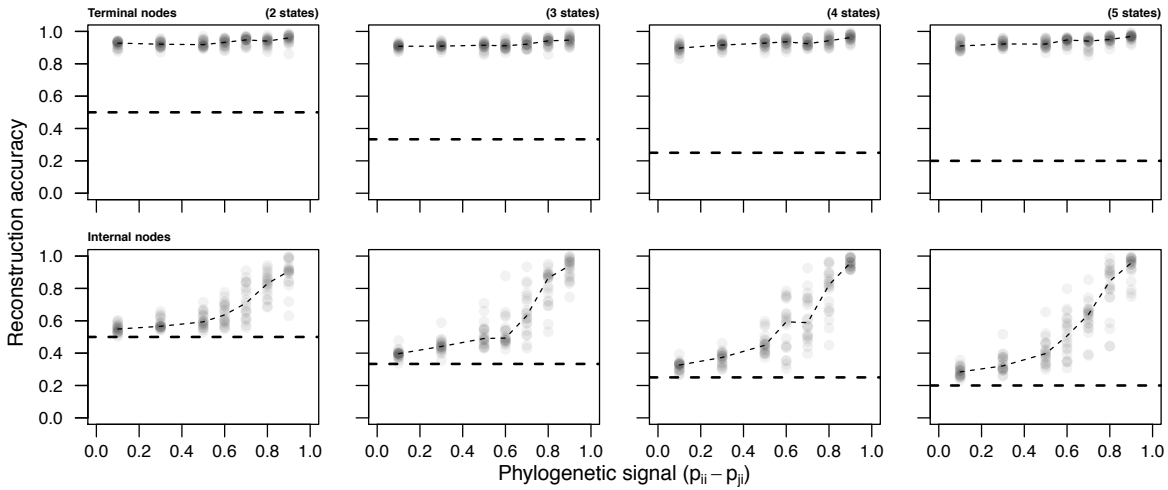
811

812

813

814

GRUNDLER AND RABOSKY



815

816 **Figure S2.** As for Figure 6 in the main text except that simulations were made using transition

817 probabilities from equation (3) rather than from equation (5).

818

819

820

821

822

823

824

825

826

827

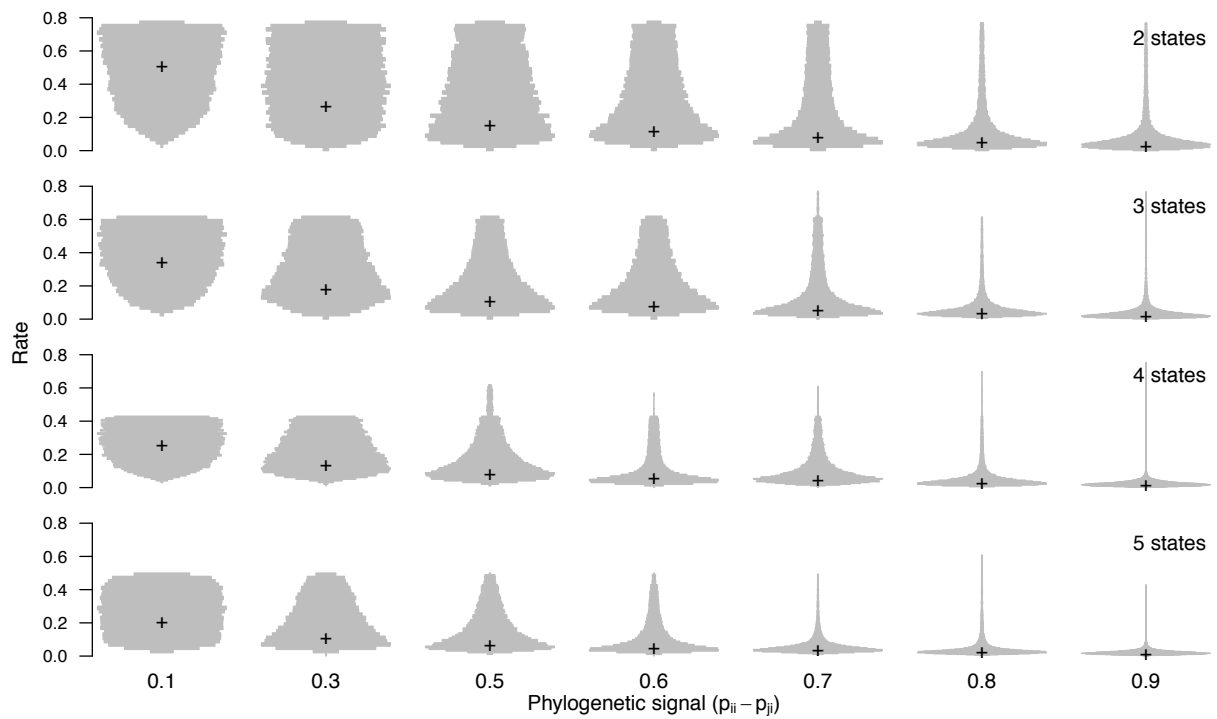
828

829

830

831

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA



832

833 **Figure S3.** As for Figure 7 in the main text except that simulations were made using transition
834 probabilities from equation (3) rather than from equation (5) and phylogenetic signal refers to the
835 phylogenetic signal of the average branch length.

836

837

838

839

840

841

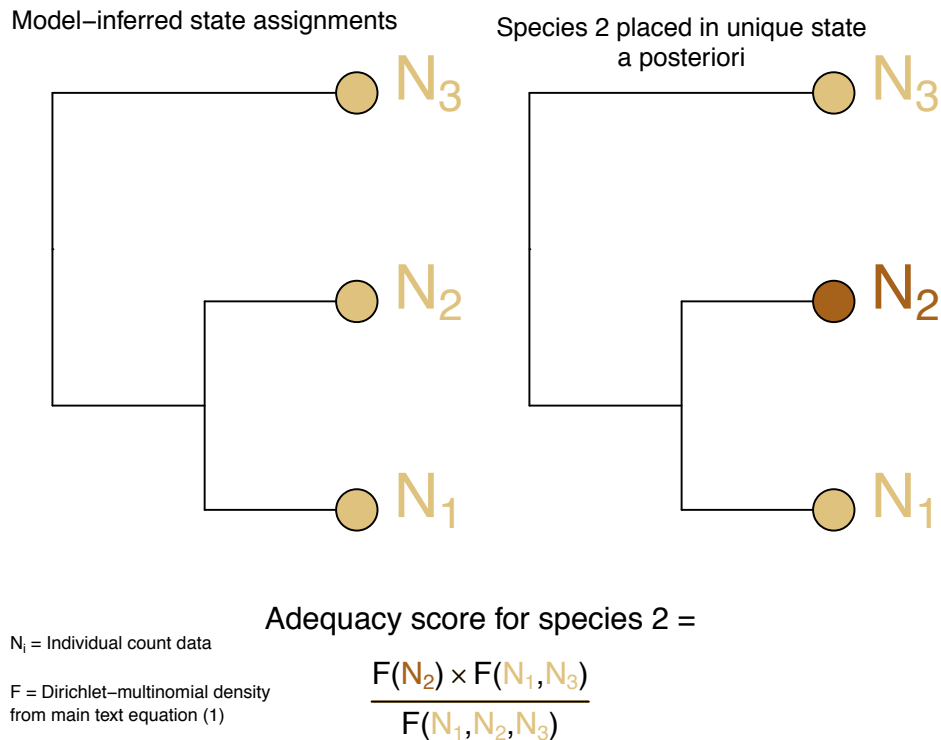
842

843

844

845

Per-species adequacy score example



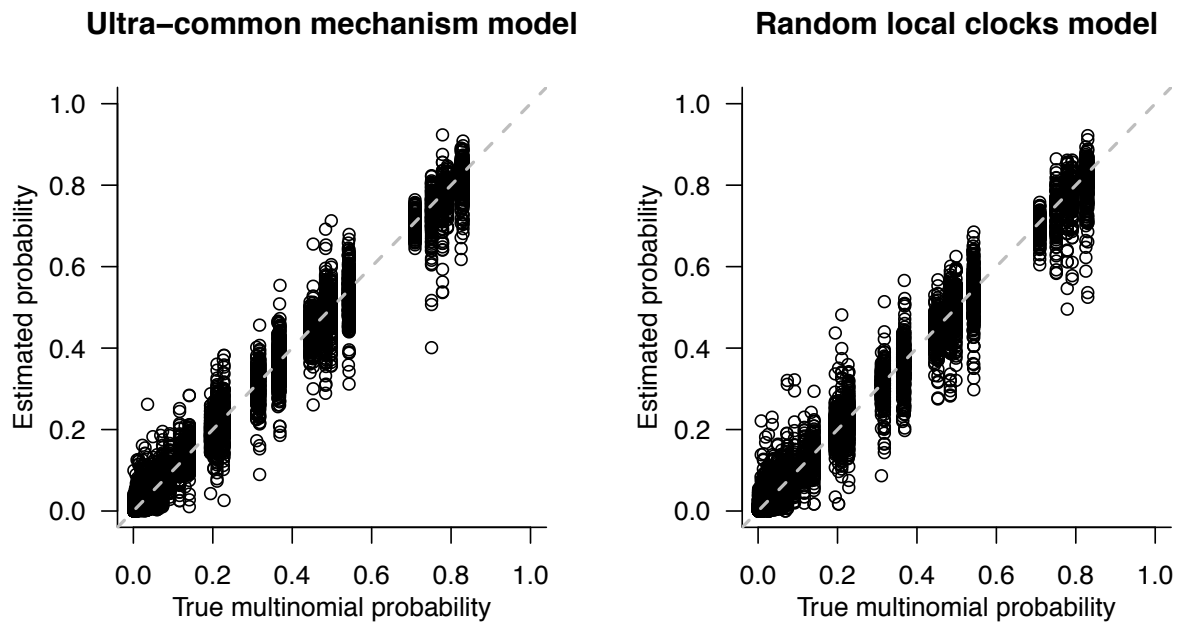
846

847 **Figure S4.** The main text defines a per-species adequacy score to assess how well the estimated
848 multinomial distributions explain sampled observations for terminal taxa. The toy example in
849 this figure depicts how such a score is calculated for species 2. In the main text, we use the
850 negative logarithm of the likelihood ratio calculation depicted in the figure. Thus, positive values
851 for the log likelihood ratio imply that the estimated multinomial is a good fit to sampled
852 observations; negative values imply the opposite.

853

854

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA



855

856 **Figure S5.** Correlation between true and estimated multinomial probabilities for datasets

857 simulated under the ultra-common mechanism model and random local clocks model.

858

859

860

861

862

863

864

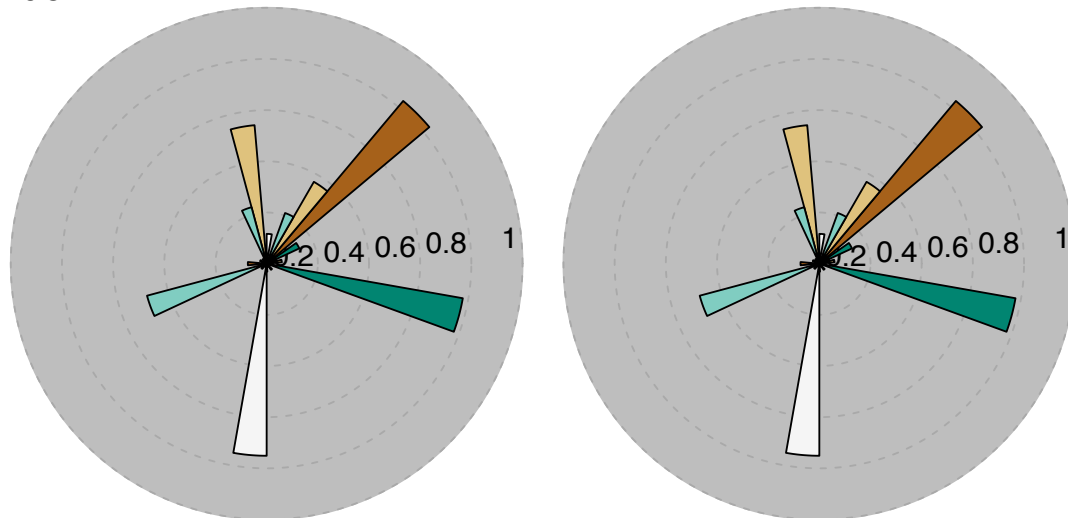
865

866

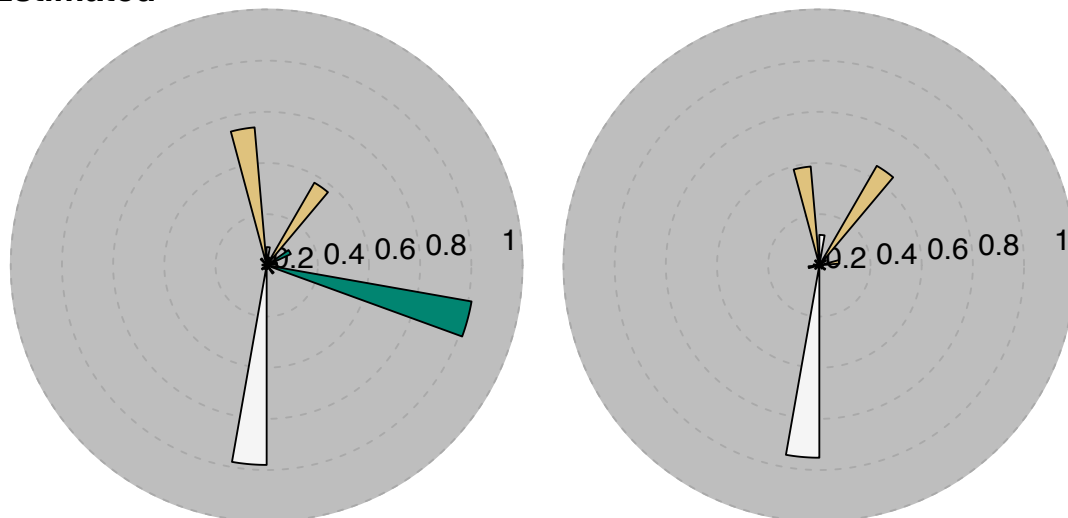
867

GRUNDLER AND RABOSKY

True



Estimated



868

869 **Figure S6.** Two examples illustrating how high accuracy of estimated multinomial proportions is

870 maintained even when the number of states is underestimated. In the left plot, the blue state only

871 generated 1 sampled observation and the brown state only generated 5 sampled observations.

872 In the right plot, the brown state only generated 8 observations; the blue state, 12; and the green

873 state, 2. In both cases, these observations are few enough in number that they do not substantially

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

874 alter the estimated multinomial proportions when they are incorporated into the count data

875 generated from the other states.

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

GRUNDLER AND RABOSKY

899 APPENDIX

900 *Derivation of $p(D|X, \beta)$*

901 We assume that conditional on the resource state the data for each species are sampled
902 independently from the latent multinomial distribution underlying the state. That is,

$$\begin{aligned} 903 \quad p(D_k|X, \theta) &\propto \prod_i p(d_i|X_i = k, \theta) \\ 904 &= \prod_i \prod_j \theta_j^{N_i^j} \end{aligned}$$

905 where we have omitted terms involving multinomial coefficients as well as the
906 dependence of θ on k for ease of notation. We also assume that for each k the multinomial
907 parameter $\theta \sim \text{Dirichlet}(\beta_1, \dots, \beta_j)$ so that,

$$\begin{aligned} 908 \quad p(D_k|X, \beta) &\propto \int_{\theta} p(D_k|X, \theta) p(\theta|\beta) d\theta \\ 909 &= \int_{\theta} \left(\prod_i \prod_j \theta_j^{N_i^j} \right) \frac{\Gamma(\sum_j \beta_j)}{\prod_j \Gamma(\beta_j)} \prod_j \theta_j^{\beta_j-1} d\theta \\ 910 &= \int_{\theta} \left(\prod_j \theta_j^{n_k^j} \right) \frac{\Gamma(\sum_j \beta_j)}{\prod_j \Gamma(\beta_j)} \prod_j \theta_j^{\beta_j-1} d\theta \\ 911 &= \int_{\theta} \frac{\Gamma(\sum_j \beta_j)}{\prod_j \Gamma(\beta_j)} \prod_j \theta_j^{n_k^j + \beta_j - 1} d\theta \\ 912 &= \frac{\Gamma(\sum_j \beta_j)}{\Gamma(n_k + \sum_j \beta_j)} \frac{\prod_j \Gamma(n_k^j + \beta_j)}{\prod_j \Gamma(\beta_j)} \int_{\theta} \frac{\Gamma(n_k + \sum_j \beta_j)}{\prod_j \Gamma(n_k^j + \beta_j)} \prod_j \theta_j^{n_k^j + \beta_j - 1} d\theta \\ 913 &= \frac{\Gamma(\sum_j \beta_j)}{\Gamma(n_k + \sum_j \beta_j)} \frac{\prod_j \Gamma(n_k^j + \beta_j)}{\prod_j \Gamma(\beta_j)} \end{aligned}$$

A HIDDEN MARKOV MODEL FOR MULTIVARIATE COUNT DATA

914 where the last equality follows because the integrand is the density function of a Dirichlet
 915 distribution with parameter $(n_k^1 + \beta_1, \dots, n_k^J + \beta_j)$. By assuming that $\beta_j = \beta$ for all j this reduces
 916 to equation (1) in the main text.

917 *Derivation of $p(X|\alpha)$*

918 Under the fully symmetric Poisson model the probability of change across an ancestral-
 919 descendant branch is,

$$920 \quad p(X_i|X_{pa(i)}, \lambda_i, t_i) = (1 - e^{-K\lambda_i t_i}) \frac{1}{K} + \delta_{X_i X_{pa(i)}} e^{-K\lambda_i t_i}$$

921 where λ_i is the branch-specific rate of evolution, t_i is the length of the branch in units of
 922 time, and $\delta_{X_i X_{pa(i)}}$ is the Kronecker delta. Because only the product $\lambda_i t_i$ matters, we can set λ_i
 923 equal to $\frac{1}{K-1}$ (meaning that the average rate is normalized to 1) so that that $t_i = v_i$ now measures
 924 time in units of expected number of changes. We assume that each $v_i \sim \text{Gamma}(\alpha, 1)$. Then,

$$\begin{aligned} 925 \quad & p(X_i|X_{pa(i)} \neq X_i, \alpha) \\ 926 \quad &= \int_0^{\infty} p(X_i|X_{pa(i)}, v) p(v|\alpha) dv \\ 927 \quad &= \int_0^{\infty} \left(1 - e^{-\frac{K}{K-1}v}\right) \frac{1}{K} \frac{1}{\Gamma(\alpha)} v^{\alpha-1} e^{-v} dv \\ 928 \quad &= \frac{1}{K} - \frac{1}{K} \int_0^{\infty} \frac{1}{\Gamma(\alpha)} v^{\alpha-1} e^{-v(\frac{K}{K-1}+1)} dv \\ 929 \quad &= \frac{1}{K} - \frac{1}{K} \left(\frac{1}{\frac{K}{K-1} + 1} \right)^{\alpha} \int_0^{\infty} \frac{\left(\frac{K}{K-1} + 1\right)^{\alpha}}{\Gamma(\alpha)} v^{\alpha-1} e^{-v(\frac{K}{K-1}+1)} dv \\ 930 \quad &= \frac{1}{K} - \frac{1}{K} \left(\frac{1}{\frac{K}{K-1} + 1} \right)^{\alpha} \end{aligned}$$

GRUNDLER AND RABOSKY

931 where the last equality follows because the integrand is the density function of a Gamma

932 distribution with parameters $(\alpha, \frac{K}{K-1} + 1)$. A similar calculation shows that

933
$$p(X_i | X_{pa(i)} = X_i, \alpha) = \frac{1}{K} + \frac{K-1}{K} \left(\frac{1}{\frac{K}{K-1} + 1} \right)^\alpha .$$

934