

## Maximum Likelihood Inference of Geographic Range Evolution by Dispersal, Local Extinction, and Cladogenesis

RICHARD H. REE<sup>1</sup> AND STEPHEN A. SMITH<sup>2</sup>

<sup>1</sup>Department of Botany, Field Museum of Natural History, 1400 South Lake Shore Drive, Chicago, Illinois 60605, USA; E-mail: rree@fieldmuseum.org

<sup>2</sup>Department of Ecology and Evolutionary Biology, Yale University, New Haven, Connecticut 06520, USA

**Abstract.**—In historical biogeography, model-based inference methods for reconstructing the evolution of geographic ranges on phylogenetic trees are poorly developed relative to the diversity of analogous methods available for inferring character evolution. We attempt to rectify this deficiency by constructing a dispersal-extinction-cladogenesis (DEC) model for geographic range evolution that specifies instantaneous transition rates between discrete states (ranges) along phylogenetic branches and apply it to estimating likelihoods of ancestral states (range inheritance scenarios) at cladogenesis events. Unlike an earlier version of this approach, the present model allows for an analytical solution to probabilities of range transitions as a function of time, enabling free parameters in the model, rates of dispersal, and local extinction to be estimated by maximum likelihood. Simulation results indicate that accurate parameter estimates may be difficult to obtain in practice but also show that ancestral range inheritance scenarios nevertheless can be correctly recovered with high success if rates of range evolution are low relative to the rate of cladogenesis. We apply the DEC model to a previously published, exemplary case study of island biogeography involving Hawaiian endemic angiosperms in *Psychotria* (Rubiaceae), showing how the DEC model can be iteratively refined from inspecting inferences of range evolution and also how geological constraints involving times of island origin may be imposed on the likelihood function. The DEC model is sufficiently similar to character models that it might serve as a gateway through which many existing comparative methods for characters could be imported into the realm of historical biogeography; moreover, it might also inspire the conceptual expansion of character models toward inclusion of evolutionary change as directly coincident, either as cause or consequence, with cladogenesis events. The DEC model is thus an incremental advance that highlights considerable potential in the nascent field of model-based historical biogeographic inference. [Ancestral state reconstruction; dispersal; extinction; Hawai'i; historical biogeography; *Psychotria*; speciation; vicariance.]

Inferring the evolution of geographic ranges of species and clades in a phylogenetic context is a major focus of historical biogeography. To this end, many questions about the past may be of interest, such as: Where were ancestors distributed? What was the tempo and direction of dispersal? How important was range expansion to lineage diversification? The breadth of potential inquiry is wide, but to date it has not been matched by an equally diverse set of methods for reconstructing range evolution on phylogenies. By and large the development of such methods has been limited to parsimony-based methods for ancestral-area optimization (Bremer, 1992, 1995; Hausdorf, 1998) and dispersal-vicariance analysis (Ronquist, 1997).

Compared to parsimony methods, development and uptake of likelihood-based approaches in historical biogeography has been slow. The contrast is particularly apparent considering how for character evolution, probabilistic models are commonly used, and the development of maximum likelihood (e.g., Schluter, 1997; Pagel, 1994, 1999) and Bayesian (e.g., Huelsenbeck and Bollback, 2001; Huelsenbeck et al., 2003; Pagel et al., 2004; Pagel and Meade, 2006) inference methods has been especially active, yielding a rich statistical tool set for investigating a wide range of macroevolutionary questions in a phylogenetic context. Quantitative statistical inference in historical biogeography has suffered from a lack of similar models and tools. It stands to reason, then, that to bring biogeographic inference up to the same level, it is worthwhile exploring analogies between characters and geographic ranges, to allow existing methods for the former guide and inform development of methods for the latter.

In character models, transitions between states are typically assumed to occur stochastically according to a Markov process, with the probability  $P_{ij}(t)$  of ancestor-descendant change from state  $i$  to state  $j$  on a phylogenetic branch of length  $t$  being a function of the model's parameter values for instantaneous transition rates. The matrix of transition probabilities for all pairs of states is generally obtained by the equation  $P(t) = e^{-Qt}$ , where  $Q$  is the instantaneous rate matrix. The likelihood function for a phylogenetic tree with observed character data at its leaves requires calculating  $P(t)$  for each branch and integrates over the conditional likelihoods of all ancestral states at every internal node, weighted by their prior probabilities (see Felsenstein, 1981, for a recursive algorithm). Transition rate parameters may thus be estimated from the tree and data by finding values that maximize the overall likelihood, without having to condition on specific character states at internal nodes. Likewise, the ancestral state of any particular node may be estimated without conditioning on other nodes in the tree, by finding the state at that node that maximizes the overall likelihood (see Pagel, 1999).

Character models have in fact already been used without modification in biogeographic studies to estimate likelihoods of ancestral areas. Nepokroeff et al. (2003) inferred historical ranges of Hawaiian *Psychotria*, understory shrubs in Rubiaceae, the coffee family, by treating the four principal island groups in the archipelago as distinct states, using models of nucleotide substitution (conveniently also having four states) to estimate ancestral areas on the phylogeny by maximum likelihood. The population terminals on their phylogeny, comprising 11 species, were all endemic to single islands,

and their model had no additional states to accommodate widespread ranges. Nepokroeff et al.'s use of maximum likelihood character state reconstruction allowed information from branch lengths to inform inferences, avoiding the potential with parsimony to underestimate evolutionary change. It also allowed uncertainty about ancestral states to be expressed in terms of probability. However, in order to properly interpret the results of such an exercise, one must decide whether models of character evolution are indeed appropriate for geographic ranges.

In character models, state transitions are not assumed to be an immediate cause or effect of speciation, and so are never inferred at internal nodes; instead, cladogenesis events are assumed to yield daughter species that initially have the same state as their parent. In the case of Nepokroeff et al. (2003), where states are geographic ranges comprising only single areas, this implies that speciation must occur within a single area, with daughter lineages identically inheriting the ancestral area (but not necessarily diverging in strict sympatry). Moreover, following a cladogenesis event, for the state of a daughter lineage to change from its parental area to another area, it must first disperse to the new area, then—instantaneously, because widespread ranges are not valid states in the model—go extinct in the original area (Ronquist, 1994). In this light, the fit of character models equating states with single geographic areas seems poor, notably in its unrealistic process of range evolution along branches, as well as in disallowing widespread ancestral ranges to evolve by subdivision at cladogenesis events (see below).

Another likelihood-based approach to the inference of ancestral geographic ranges was described by Ree et al. (2005), which departed from previous methods in two significant ways. First, they modeled geographic range evolution by stochastic dispersal and local extinction events in a set of discrete areas in continuous time, enabling for the first time the calculation of  $P_{ij}(t)$  for phylogenetic tree branches where  $i$  and  $j$  represent ancestor-descendant geographic ranges. Unlike Nepokroeff et al. (2003), widespread ranges encompassing two or more areas are valid states for single species in their model, being the direct outcome of dispersal events. Second, with explicit reference to the spatial consequences of speciation, they considered distinct scenarios of range inheritance at cladogenesis events to be the fundamental “ancestral states” of interest at internal nodes on phylogenetic trees. In particular, they reasoned that for a widespread ancestor, lineage divergence could occur either between a single area and the remainder of its range, or within an area, e.g., by peripheral isolate speciation, and described a flat prior expectation for each kind of divergence. Range evolution as a direct consequence of speciation contrasts directly with character models, which assume that states do not change coincidentally with cladogenesis events. Moreover, divergence within an area allows persistence of a widespread ancestral range through a speciation event (see fig. 3 in Ree et al., 2005), highlighting how their pool of range inheritance scenarios differs from

dispersal-vicariance analysis (Ronquist, 1997) in allowing more ways for widespread ranges to be subdivided and inherited at nodes than solely the vicariance pattern of divergence between areas.

Ree et al.'s method was novel because, in describing an explicit model of dispersal and local extinction underlying transition probabilities  $P(t)$  for geographic ranges and assigning prior probabilities to ancestral states (range inheritance scenarios) at nodes, it provided a means of calculating the biogeographic likelihood of a phylogenetic tree with range data arrayed at its tips in a manner directly analogous to character states. However, it differed from character models in one important way: it was prohibitively slow for practical applications of maximum likelihood optimization to estimate either model parameters (rates of dispersal and local extinction) or ancestral range inheritance scenarios, because it relied critically on a time-consuming simulation step in computing  $P(t)$ . In developing the method, Ree et al. emphasized the importance of bringing a variety of temporal information sources other than just the tree and extant range data to bear on historical inferences, including fossils, sea levels, climate, continental positions, etc., in effect, any data relevant to where the organisms of interest may have occurred. Because such external information was envisioned to yield arbitrarily complex functions for dispersal and extinction rates through time, general analytical solutions for  $P(t)$  were forgone in favor of a Monte Carlo approach, in which relative frequencies of simulation outcomes were used to estimate range transition probabilities on tree branches. A major advantage of the simulation approach lies in its flexibility: any number of temporal or spatial constraints on range evolution may be imposed, either by building them into the simulation engine or using them as filters for accepting valid outcomes. Ree et al. (2005) demonstrated how temporal information on land bridges in the Northern Hemisphere could be incorporated into simulations estimating  $P(t)$ , thus influencing likelihood calculations and ancestral range inferences. However, the computational burden of simulation severely limits the general applicability of their method.

In this paper, we describe how to refine the method of Ree et al. (2005) to overcome this drawback, showing how instantaneous rates of dispersal and local extinction can be used to construct a rate matrix for transitions between geographic ranges that is directly analogous to the rate matrix  $Q$  typically used in character models. This offers an analytical solution to ancestor-descendant transition probabilities  $P(t)$ , but likelihoods under the original and new methods are otherwise calculated identically. The new method makes feasible parameter estimation and ancestral state reconstruction by maximum likelihood optimization using commonly used methods developed for characters, and we test the accuracy of these inferences against simulated data. Substantial flexibility is retained in allowing conditions to be imposed on the range evolution process, as we demonstrate by modeling changes in dispersal opportunity through time. As an illustrative empirical case study,

we revisit and apply the method to Nepokroeff et al.'s (2003) study of Hawaiian *Psychotria*.

#### THE DISPERSAL-EXTINCTION-CLADOGENESIS (DEC) MODEL

In this section we describe the stochastic model of geographic range evolution, referred to hereafter as the dispersal-extinction-cladogenesis (DEC) model, emphasizing how it differs from that described by Ree et al. (2005) in the calculation of transition probabilities from an instantaneous rate matrix. Readers should consult the former paper for the original description of model features retained identically here, particularly in how range inheritance scenarios are enumerated and assigned prior probabilities.

*Geographic ranges.*—We represent the geographic range of a species as a string denoting its presence in a set of areas. For three areas labeled 1, 2, and 3, the set of possible ranges is thus  $\emptyset, 1, 2, 3, 12, 13, 23, 123$ . With the exception of the empty range ( $\emptyset$ ), these comprise all theoretically observable states for extant species.

*Range evolution along phylogenetic branches.*—In the absence of lineage divergence, the range of a species evolves by two stochastic processes: dispersal between areas (range expansion) and local extinction within areas (range contraction). Parameters for these are  $D_{ij}$ , the rate of dispersal from area  $i$  to area  $j$ , and  $E_i$ , the rate of local extinction in area  $i$ . For  $n$  areas, the most general model would have  $n^2 - n$  independent (free) parameters for dispersal rates: one for each direction between each pair of areas. In addition, it would have  $n$  free parameters representing area-specific local extinction rates. Simpler models would assume fewer free parameters, with the simplest having a single dispersal rate and a single local extinction rate, each uniform across the areas in the model and across the phylogeny.

These rate parameters can be used to construct the matrix of instantaneous transition rates between geographic ranges ( $Q$ ). For three areas,

$$Q = \begin{bmatrix} & \emptyset & 1 & 2 & 3 & 12 & 13 & 23 & 123 \\ \emptyset & - & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & E_1 & - & 0 & 0 & D_{12} & D_{13} & 0 & 0 \\ 2 & E_2 & 0 & - & 0 & D_{21} & 0 & D_{23} & 0 \\ 3 & E_3 & 0 & 0 & - & 0 & D_{31} & D_{32} & 0 \\ 12 & 0 & E_2 & E_1 & 0 & - & 0 & 0 & D_{13} + D_{23} \\ 13 & 0 & E_3 & 0 & E_1 & 0 & - & 0 & D_{12} + D_{32} \\ 23 & 0 & 0 & E_3 & E_2 & 0 & 0 & - & D_{21} + D_{31} \\ 123 & 0 & 0 & 0 & 0 & E_3 & E_2 & E_1 & - \end{bmatrix}. \quad (1)$$

(For clarity, geographic ranges at the start and end of a transition are shown in order of increasing size along the first column and row, respectively.) We assume that over an infinitesimal time interval, only one event may occur, so transition rates between ranges that differ by more than one dispersal or extinction event are set to zero.

Non-zero cells in the matrix below the diagonal represent range contractions by local extinction; those above the diagonal represent range expansions by dispersal. Some range expansions involve a sum of rates; e.g., the transition from range 12 to 123 involves dispersal rates from areas 1 and 2 to area 3. In general, the rate of expansion from a starting range  $r$  to a wider range  $r'$  is the sum of dispersal rates from all areas in  $r$  to the area of expansion in  $r'$ . The rate of contraction from  $r'$  back to  $r$  is the rate of local extinction in that area. The elements along the diagonal of the rate matrix are defined such that the sum of rows is equal to zero.

Constructing the rate matrix  $Q$  algorithmically in this manner brings us back to the equation for range transition probabilities as a function of time,  $P(t) = e^{-Qt}$ , which can be computed much faster (by orders of magnitude) than the simulation-based method of Ree et al. (2005). However, recall that in order to calculate the likelihood of the data given the assumed tree, additional information is needed; namely, an enumeration of possible ancestral states at internal nodes and their prior probabilities.

*Range evolution at cladogenesis events.*—Following Ree et al. (2005) and as described in the introductory section, we assume that if an ancestor is widespread across two or more areas, speciation can happen in one of two ways: lineage divergence could arise either between a single area and the rest of the range, or within a single area. This leads to nonidentical range inheritance, with one daughter species always inheriting a single-area range, and the other inheriting either the remainder of the ancestral range or its entirety, respectively. Ree et al. (2005) described how a flat prior for the ancestral range can be multiplied by a flat prior for between- and within-area divergence patterns to obtain the overall prior for each range inheritance scenario.

It is worth noting that the above assumes cladogenesis events that are strictly bifurcating; i.e., that speciation gives rise to two, and only two, descendant lineages. We treat this as the general case and view simultaneous divergence of more than two species from a common ancestor (as depicted by “hard” phylogenetic polytomies) as exceptional. For the DEC model to include such cases, an explicit rationale for enumerating polytomous range inheritance scenarios would be needed. By contrast, “soft” polytomies indicative of phylogenetic uncertainty do not require such extensions of the model and may be accounted for statistically (see Discussion).

*Inferring ancestral ranges by maximum likelihood.*—With a matrix of range transition rates ( $Q$ ) derived from rates of dispersal between areas and rates of local extinction within areas, and prior probabilities for range inheritance scenarios, we now have all the components necessary to calculate the likelihood of a phylogenetic tree with observed range data arrayed at its tips. This is done exactly as for character data but integrating over the conditional likelihoods of range inheritance scenarios rather than ancestral character states at internal nodes. The close correspondence of the DEC model to character models allows us in principle to use all maximum likelihood methods that are normally applied to character inference. In the

following sections we focus specifically on estimating model parameters (rates of dispersal and local extinction) and ancestral range inheritance scenarios, following Schluter et al. (1997) and Pagel (1999).

### *Testing Inferences Against Simulated Data*

**Methods.**—To explore inferences about range evolution in a controlled setting, we programmed a simulator for concurrent range evolution and phylogenetic tree growth, combining the DEC model with stochastic cladogenesis. We used this biogeographic birth-death model to grow phylogenetic trees with known histories of dispersal and local extinction. The model had three areas, with a single rate parameter for both dispersal between areas and local extinction within areas (dispersal and extinction were thus constrained to be equal). Across simulations, this rate was allowed to vary between 0.01 and 0.2, whereas the speciation rate was held constant at 0.4. Branch lengths on simulated trees are thus in units of expected dispersal and local extinction events. Between-area and within-area subdivision of widespread ancestors were equiprobable at cladogenesis events. For each simulation the geographic range at the root node was randomly initialized to a single area, and the tree was allowed to grow to a predefined size (number of extant leaf nodes).

Two thousand trees were simulated in each of two sizes, 20 and 100 extant leaf nodes. To approximate empirical studies of extant species, trees were pruned of branches that went globally extinct as a result of stochastic local extinction. For each pruned tree, optimal dispersal and extinction rates were first estimated by global maximum likelihood, and these values were then used to estimate most likely range inheritance scenarios at internal nodes. Scenarios at each node were considered in isolation, without conditioning on scenarios at other nodes.

**Results.**—Estimated rates for dispersal and local extinction are shown compared to their apparent rates in Figure 1. By apparent rate we mean the number of simulated events that actually occurred on the pruned tree, divided by the tree length. Both dispersal and local extinction tend to be underestimated, with estimated dispersal appearing to be lower than apparent by a constant proportion, and extinction rates being rarely estimated far from zero across the range of apparent rates. Additional studies are needed to fully understand these patterns, and a more in-depth analysis of simulation results is in preparation.

Despite inexact rate estimates, ancestral range inheritance scenarios can be reconstructed with considerable success if dispersal and local extinction events are rare relative to speciation. Accuracy declines as the overall rate of range evolution increases: for trees with estimated values of dispersal plus extinction spanning 0 to 0.05, binned into five 0.01-unit intervals, the mean proportion of nodes reconstructed accurately per 20-taxon tree declines from 0.957 in the first bin to 0.724 in the last,

with similar numbers for 100-taxon trees (Fig. 2). We emphasize that this relationship between accuracy and rate can be interpreted properly only in light of the value for speciation rate, 0.4, which was held constant across all simulations. Further work is needed to be able to generally predict type I and II error rates in empirical studies.

### *Incorporating Prior Conditions and Constraints*

We return now to the point made by Ree et al. (2005) that many factors, both biotic (e.g., dispersal traits) and abiotic (e.g., distances between areas), may influence our prior expectations about range evolution. They emphasized the potential for incorporating external information into biogeographic models, especially with regard to their effect on dispersal opportunities through time (e.g., windows of land bridge availability). Probabilities for dispersal success can, at least in theory, be specified as arbitrarily complex, continuous-time functions encapsulating any and all factors deemed relevant to the dispersal history of the clade of interest. It is straightforward to extend this approach to other aspects of the model; e.g., to develop functions for area- or lineage-specific extinction probabilities. Such functions are readily incorporated into the previously described simulation-based method for estimating  $P$ .

For the DEC model, incorporating prior information on range evolution requires a somewhat different strategy. Variation in dispersal and extinction can be specified by time periods in which particular constraints and conditions apply to parameter values. For example, the time of origin  $t_0$  of a volcanic island could be used to demarcate two periods corresponding to its availability for colonization: in the earlier period, all dispersal rates to the island are set to zero, and the set of valid ranges in the model is reduced to those that exclude the island. The later period is unconstrained. The periods thus stratify the phylogeny of interest, dividing it into branches (and segments of branches) that fall on either side of  $t_0$ . Calculating the likelihood of extant range data on a phylogeny spanning these strata is then done by iteratively working backwards in time. Starting with the most recent period and its unconstrained model of range evolution, fractional likelihoods for ranges are evaluated at points along branches intersecting  $t_0$ , conditional on those ranges being valid outcomes of the preceding period (i.e., ranges excluding the island originating at  $t_0$ ). These fractional likelihoods are then input into likelihood calculations for parts of the tree prior to  $t_0$ , using probabilities of range transitions from the constrained model. In general, to evaluate the likelihood of extant range data on a tree that is stratified according to period-specific constraints, all that is required is a postorder (tips-to-root) traversal of the tree that, in addition to evaluating fractional likelihoods of range inheritance scenarios at lineage divergence points, also evaluates fractional likelihoods of ranges at points where branches cross period boundaries.

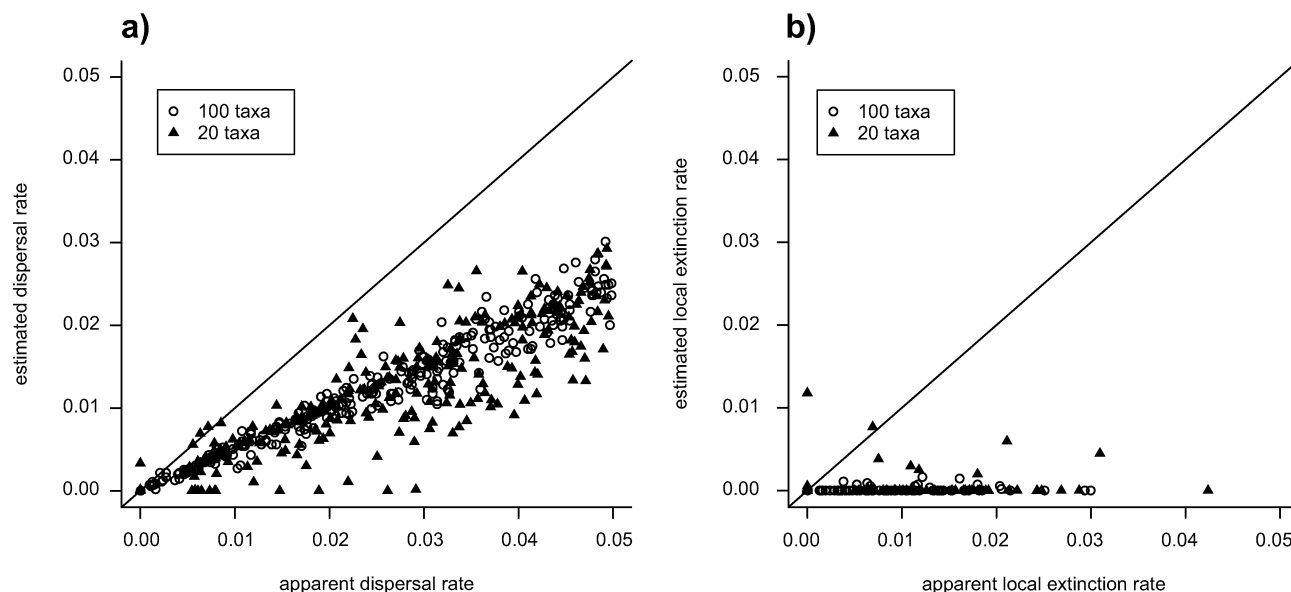


FIGURE 1. Maximum likelihood parameter estimates for rates of dispersal (a) and local extinction (b) under the DEC model from trees of two size categories simulated according to a geographic birth-death model (see text), then pruned of extinct branches. For each size, a random sample of 300 out of 2000 trees are shown. Apparent rates (the number of actual events simulated on a pruned tree divided by tree length) are plotted against estimated rates, with points below the diagonal representing underestimates. Bias toward underestimation is evident for both dispersal and local extinction, with estimates of the latter consistently being close to zero.

### Implementation

The DEC model and routines for estimating dispersal and extinction parameters and ancestral range inheritance scenarios are available in the software package *lagrange*, distributed by the authors from <http://code.google.com/p/lagrange>.

### RANGE EVOLUTION IN HAWAIIAN PSYCHOTRIA

#### Objectives

We revisited the case of Hawaiian *Psychotria* (Nepokroeff et al., 2003) as an empirical opportunity to apply the current model, using the four areas defined in that study: (1) Kaua'i, including Ni'ihau; (2) O'ahu; (3) Maui Nui (including Moloka'i, Lana'i, Maui, and Kaho'olawe); and (4) Hawai'i (Fig. 3a), which we label K, O, M, and H. Our first objective was to fit their data to a general four-area DEC model and explore inferences on the phylogeny without imposing any temporal conditions on range evolution corresponding to times of island origin. Specifically, under this model we wished to determine (1) which area was most likely colonized by the ancestor of the clade; (2) the relative degree of within-area versus between-area lineage divergence following dispersal, as inferred from ancestral range inheritance scenarios (i.e., the extent to which dispersal leads to cladogenesis); and (3) whether simpler models, e.g., those allowing only dispersal between adjacent islands, were favored over the general model. Our second objective was to construct a temporally stratified model of range evolution reflecting absolute ages of the Hawaiian islands and then compare inferences between models with and without such temporal constraints.

### Methods

We used the maximum likelihood, ultrametric molecular phylogeny from Nepokroeff et al. (2003), which included 22 in-group populations representing the 11 recognized endemic species. We modified it by scaling its length to absolute time and editing the topology to accommodate the model's requirement that the tree be completely bifurcating. On the assumption that ancestral colonization occurred soon after the origin of Kaua'i, the oldest extant Hawaiian island, we set the root age of the tree at 5.1 Myr. To resolve polytomies, we combined the two O'ahu populations of *P. mariniana* into a single terminal and combined the three Maui Nui populations of *P. mariniana* into a single terminal, placing the latter as sister to the Hawai'i population of *P. hawaiiensis*. We also grouped the Maui Nui populations of *P. kaduana* and *P. mauiensis* as sister taxa and assigned a very short length ( $10^{-5}$ ) to the subtending branch of that clade. These modifications are somewhat arbitrary (e.g., the ancestor could have colonized an island now submerged, earlier than 5.1 Ma), so we emphasize that the empirical results of our re-analysis should be treated with appropriate caution, with more weight given to the illustrative value of the exercise.

We constructed a stratified model of four discrete time periods corresponding to approximate maximum ages of Kaua'i (5.1 Myr), O'ahu (3.7 Myr), Maui Nui (1.9 Myr), and Hawai'i (0.5 Myr; Carson et al., 1995). In each period, geographic ranges and dispersal between areas were restricted to include only extant islands but were otherwise unconstrained; across periods, rates of dispersal and local extinction were assumed to be constant.

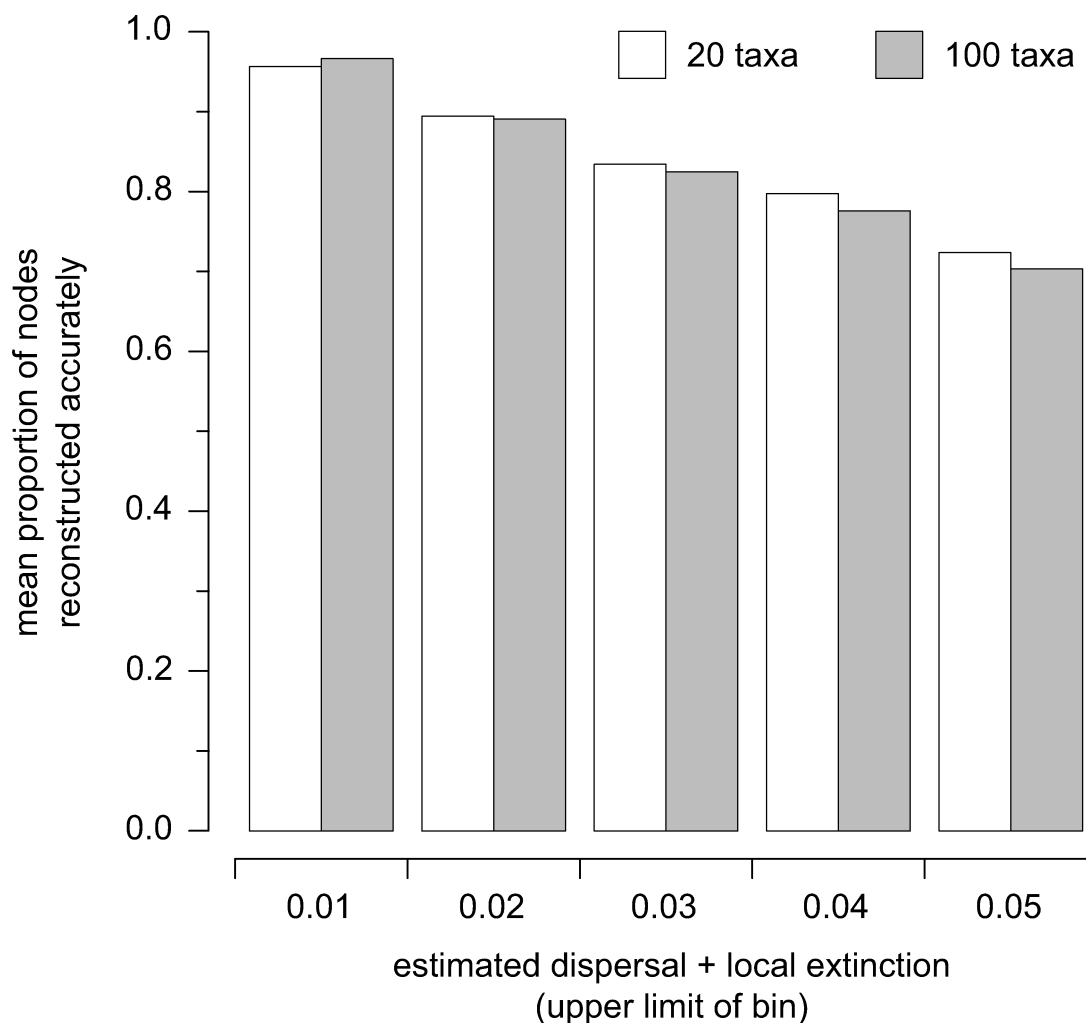


FIGURE 2. Proportion of nodes accurately reconstructed by maximum likelihood for ancestral range inheritance scenarios at cladogenesis events on simulated trees, in relation to the sum of estimated rates of dispersal and local extinction. Rates are binned in intervals of 0.01, with sample size (number of simulated trees) per bin for each tree size ranging from 246 to 574. All trees were simulated with a speciation rate of 0.4.

At the root of the tree, we considered only single-area ranges, reasoning that colonization originally occurred on a single island. The likelihood of each area at the root was estimated using locally optimal rates of dispersal and extinction (i.e., rates were optimized to their maximum likelihood values conditional on the root state), corresponding to the “local” method of estimating ancestral character states by maximum likelihood (see Pagel, 1999). The optimal root area and rates were then fixed for comparing likelihoods of alternative range inheritance scenarios at internal nodes further up the tree. At each node, the likelihood of each scenario was calculated without conditioning on scenarios at any other nodes in the tree.

### Results

In the unconstrained model, Kaua’i is the most likely island of colonization by *Psychotria*, with other areas having successively lower likelihoods (Table 1), all outside the confidence window of two log-likelihood

units (Edwards, 1992). Similarly, at most internal nodes, optimal range inheritance scenarios also score significantly better than any alternatives, but in some cases other scenarios are also statistically plausible, indicating localized uncertainty. With this in mind, for simplicity we show and discuss only the optimal reconstruction (Fig. 3b), which reveals five branches along which ranges evolve, all involving dispersal events from single-area ancestral ranges. Two of the branches also involve local extinction events. The predominant pattern is that widespread ranges do not persist for long before being split by allopatric cladogenesis or reduced by local extinction. Three of the branches terminate in between-area (vicariant) lineage divergence, suggesting that dispersal was a proximate cause of speciation in those cases or, in other words, that colonization of a new island led to rapid coalescence of a distinct lineage. By contrast, one instance of dispersal (O to OM, along the branch subtending the *P. kaduana* clade) leads to within-area lineage divergence, with the widespread range OM

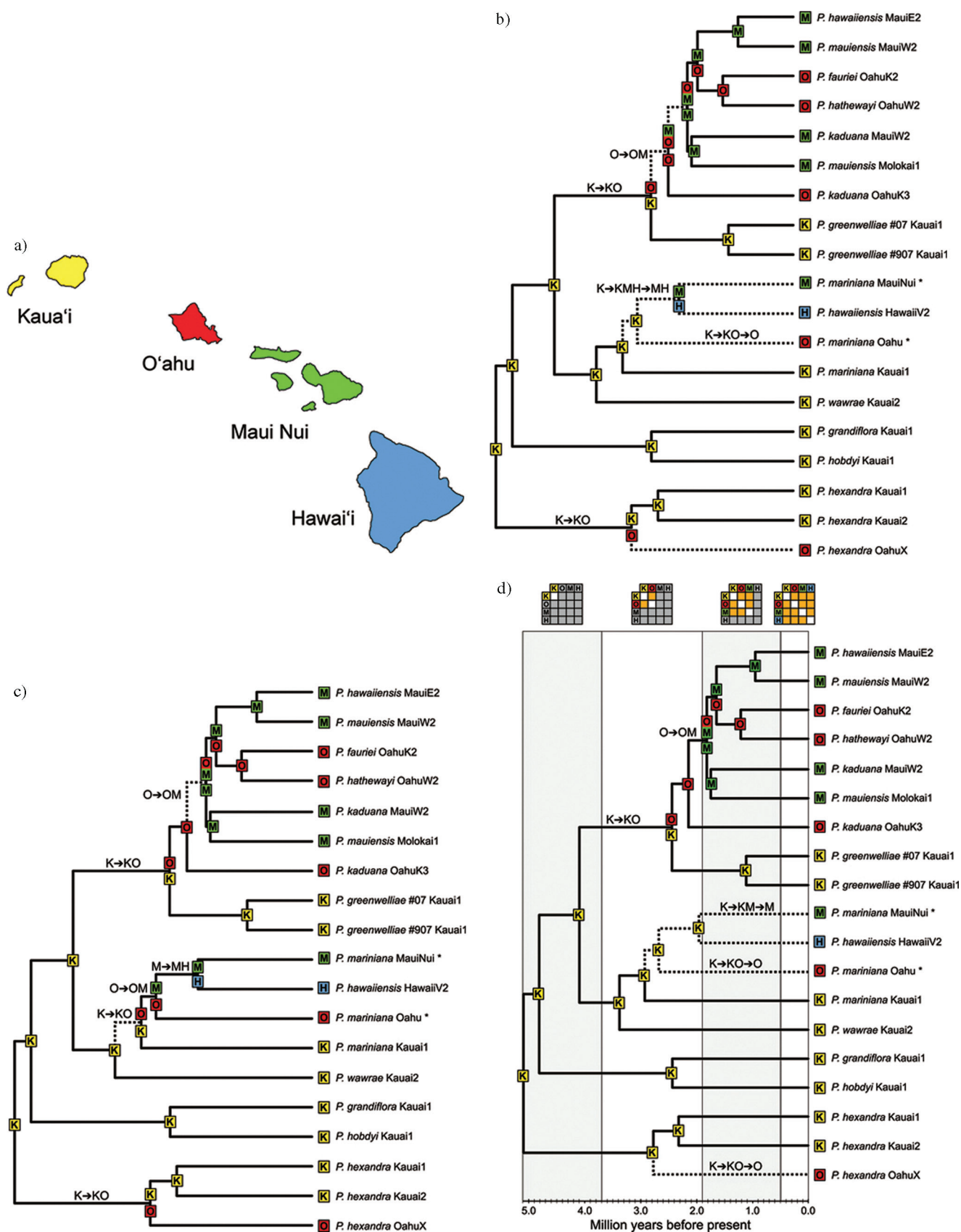


FIGURE 3.

TABLE 1. Inferences about the ancestral area and range evolution parameters of Hawaiian *Psychotria* under DEC models. The unconstrained model (M0) allows geographic ranges to include any combination of islands in the archipelago and permits direct dispersal between any pair of islands. M1 and M2 restrict ranges to include a maximum of two adjacent islands. M2 further limits dispersal to be eastward between adjacent islands. The stratified model permits dispersal to islands only after their time of geological origin, thus with a root age of 5.1 Ma, the only ancestral area possible is Kaua'i.

Model	Area	−ln(L)	Dispersal	Extinction
M0	Kaua'i	35.758	0.040	0.0358
	O'ahu	40.700	0.041	0.024
	Maui Nui	44.378	0.054	0.076
	Hawai'i	45.323	0.058	0.085
M1	Kaua'i	34.636	0.093	0.017
	O'ahu	38.877	0.112	0.052
	Maui Nui	48.683	0.207	0.164
	Hawai'i	55.396	0.377	0.280
M2	Kaua'i	32.434	0.132	0.009
	O'ahu	106.018	0.174	0.103
	Maui Nui	107.701	0.216	0.101
	Hawai'i	118.930	0.173	0.066
Stratified	Kaua'i	40.777	0.075	0.082

being retained after speciation. In this case, both of the branches that retain OM are relatively short, terminating in between-area divergence at the base of the clade containing *P. hathewayi*, further suggesting that range expansion leads to rapid lineage divergence in allopatry.

Two other notable patterns emerge from the reconstruction of biogeographic history shown in Fig. 3b. First, four out of the five inferred dispersal events involve areas adjacent to each other in the Hawaiian chain and are consistently eastward from older to younger islands. Second, widespread ancestral ranges generally include only two adjacent areas and, as noted above, apparently decay rapidly to single areas. These patterns suggest that models more constrained than the initial one, in terms of dispersal opportunity and allowed ranges that correspond to the spatial layout of the Hawaiian chain, may fit the data better.

To investigate this by a posteriori exploration, we reestimated the likelihood of the data under two alternative models (M1 and M2) for comparison to the original model (M0). Both M1 and M2 constrained ancestral ranges to a maximum of two adjacent areas (the set  $\{\emptyset, K, O, M, H, KO, OM, MH\}$ ), effectively having a new parameter, fixed at zero, for the prior probabili-

ties of other ranges. M2 constrained dispersal to be only eastward to adjacent areas, effectively having another parameter, also fixed at zero, for dispersal rates westward and between nonadjacent islands. These models are not nested in terms of free parameters, as they are the same in having one for dispersal and one for local extinction; rather, they differ in having additional parameters fixed at a boundary value. So although we were guided by initial inferences in constructing M1 and M2, there is no a priori expectation for one of the models to yield higher or lower likelihood scores than the others.

Log-likelihood scores of the data under M1 and M2 were both better than under M0 (Table 1), with M2 being the best by 2.21 log-likelihood units. Moreover, inferences of range evolution under M2 exhibit less uncertainty compared to M0 (Fig. 3c). Thus, modifying the initial model by inspecting unconstrained inferences leads to a more refined model that better fits the data, independently confirming expectations of "progression rule" dispersal based on the spatial arrangement and relative ages of the Hawaiian islands (Carson and Clagne, 1995; Funk and Wagner, 1995).

*Island ages and dispersal constraints.*—Under the stratified model that disallows dispersal to islands until their time of geological origin, the likelihood is substantially lower relative to the others (Table 1), but to the extent that phylogenetic divergence times and island ages are accurate, inferences of range evolution under this model may be more realistic than unconstrained inferences; we return to this issue in Discussion. Maximum likelihood reconstructions of range inheritance scenarios are depicted in Fig. 3d. A notable difference is at the base of the *P. kaduana* clade: unstratified models yielded some uncertainty in whether the ancestral range included Maui Nui, O'ahu, or both; here, the node occurs prior to Maui Nui, resulting in significant confidence that the ancestral area was O'ahu only. Nepokroeff et al. (2003) also found ambiguous inferences here (node 36 in their fig. 4) with both parsimony and likelihood methods.

## DISCUSSION

### *Methodological Links between Historical Biogeography and Character Evolution*

In this paper we attempt to bridge a conceptual gap in phylogenetic biology between models and inference

FIGURE 3. (a) Map of the Hawaiian archipelago showing the four areas (islands or island groups) of interest in the reanalysis of *Psychotria* biogeography: Kaua'i, O'ahu, Maui Nui, and Hawai'i. (b) Maximum likelihood reconstruction of geographic range evolution under the unconstrained DEC model. (c) Maximum likelihood reconstruction of geographic range evolution under the simplified DEC model (M2; see text) allowing only ranges with one area (K, O, M, H), or two areas adjacent along the Hawaiian chain (KO, OM, MH), with directional, west-to-east "stepping-stone" dispersal. Relative to the unconstrained model, M2 yields a significantly higher likelihood score, with less uncertainty about ancestral range inheritance. (d) Maximum likelihood reconstruction of geographic range evolution under the constrained DEC model disallowing colonization of areas before their origin. Above each of the four time periods, possible area-to-area dispersal is depicted as a matrix, with orange and gray cells indicating yes and no, respectively. In b to d, ranges inherited from widespread ancestors following cladogenesis events are shown at the bases of diverging branches; single-area ancestral ranges are shown at nodes. Dashed-line branches indicate those for which alternative reconstructions at the base fall within two log-likelihood units of the optimal scenario shown. Range transitions along branches show the inferred sequence of dispersal and extinction events. Taxa with names followed by asterisks represent combinations of taxa from Nepokroeff et al. (2003; see text).



methods for character evolution and those for geographic range evolution. We show that constructing an instantaneous rate matrix of transitions between ranges from rates of dispersal and local extinction yields a model directly analogous to evolutionary models for discrete character states: namely, one describing rates and probabilities of transitions between states (ranges) along phylogenetic branches, applicable to maximum likelihood inference of states (range inheritance scenarios) at cladogenesis events. Creating this link between characters and ranges should pave the way for a variety of comparative methods to be ported from the realm of discrete trait evolution to that of historical biogeography. Here we demonstrate simple inference of ancestral states at individual nodes on a given tree by maximum likelihood, but this is clearly a minute fraction of the range of possibilities. For example, hypotheses about directional trends in dispersal could be easily examined using likelihood ratio tests, by comparing models differing in the configuration of free parameters in the dispersal rate matrix (Pagel, 1999). It should be similarly straightforward to use the DEC model in a Bayesian context; e.g., with the use of stochastic mapping (Nielsen, 2002; Huelsenbeck et al., 2003), a method of statistically sampling complete sequences of evolutionary change that is particularly well-suited to applications that integrate over phylogenetic uncertainty in topology and branch lengths. Applying stochastic mapping to range evolution would allow statistical inferences to be drawn from the posterior probability distribution of biogeographic histories (dispersal, extinction, and range inheritance events), without conditioning on point estimates of phylogeny or evolutionary rates. The extent to which biogeographic inferences under the DEC model are robust to uncertainty in these items (what might be considered nuisance parameters for some hypothesis tests) is highly data-dependent, and so historical biogeographers should follow the ample precedent set in the literature on character evolution to account for this uncertainty.

The potential benefits of the DEC model may also flow in the opposite direction, from range evolution to character evolution. For example, the DEC model could serve as a basis for incorporating evolution at cladogenesis events into character models, which are currently uniform in assuming that lineage divergence occurs with identical inheritance of the ancestral condition. The analogy of range evolution highlights the potential for characters to evolve as an immediate effect of speciation, particularly if character polymorphism or variance is viewed as equivalent to being geographically widespread.

The effect of ranges and characters on the rate of cladogenesis itself is another factor of common interest in need of theoretical development. For characters, phylogenetic models have by and large only attempted to correlate particular states with distinct rates of diversification (e.g., Slowinski and Guyer, 1993; Paradis, 2005; Ree, 2005); state change as a direct cause of cladogenesis events remains unexplored territory. In biogeography, it stands to reason that as ranges expand, the probability of allopatric divergence should increase; however, this

intuition has yet to be incorporated into stochastic models for phylogenetic inference. New methods involving numerical integration of differential equations are likely to prove useful in this area (Maddison et al., 2007).

### *Constraints on Range Evolution*

Historical biogeography is grounded in the notion that Earth history has profoundly influenced the geographic ranges of species or, in other words, that the evolutionary histories of areas and lineages are tightly coupled. This is the basic rationale of vicariance biogeography, and it is also the reason why Ree et al. (2005) emphasized the importance of bringing information on Earth history, particularly regarding area connections through time, to bear on inferences of range evolution. Here we have shown how the DEC model can incorporate such information in applying time-dependent constraints on dispersal. The specific case, origins of volcanic islands dictating when colonization could have happened, yields time periods during which dispersal constraints affect all extant lineages equally. Ree et al. (2005) suggest several other kinds of constraints that could be applied, not just across time periods but specific to particular lineages (e.g., corresponding to differences in traits conferring dispersal ability and hence rate) or areas (e.g., size or habitat considerations influencing local extinction rates).

Such constraints generated from external information are easily imposed, but whether they are actually supported by the comparative data at hand is another matter. In general, it would seem prudent at the outset to avoid unduly weighting one kind of evidence over another. In some cases, a poor fit of comparative data to a constrained model might simply be due to bad point estimates; for example, we constrained *Psychotria* range evolution using fixed dates for island origins and lineage divergences, resulting in a significant decline in likelihood compared to the unconstrained model. A search for higher likelihoods within the bounds of error in those estimates could conceivably yield more realistic reconstructions of biogeographic history. In other cases, constraints might best be posed as hypotheses to be tested, e.g., whether winged fruits confer higher dispersal rates. In the context of the DEC model, this could be done using likelihood ratios, comparing models with constant dispersal rates to those with rates uncoupled between winged and nonwinged clades.

Moving even further away from the end of the spectrum where geological information trumps phylogenetic information, one might imagine optimizing likelihoods of unconstrained range evolution across multiple clades inhabiting the same set of areas to discover whether phylogenetic signatures of Earth history are indeed evident; i.e., whether geographic constraints themselves can be discovered from comparative data. For example, the temporal boundaries of a land bridge (its window of dispersal opportunity) could become free parameters to be optimized by maximum likelihood using a DEC model applied to multiple clades.

### Practical Considerations in Empirical Studies

The DEC model is more complex than character models in that “states,” whether distinct ranges or inheritance scenarios, are more numerous than the basic units of interest, namely areas. The size of the range transition rate matrix increases exponentially with the number of areas in the model if all presence-absence combinations are allowed; for large rate matrices, the likelihood calculation may incur a substantial computational burden, hindering optimization. For empirical work it is thus advantageous to avoid dissecting the Earth too finely. Circumscribing areas in practice should balance biotic and abiotic factors, consider area identities over the time period of interest, and aim for geographic granularity appropriate to the inference objectives. It may also be worthwhile to consider reducing the set of allowed ranges. In this we recommend taking an empirical approach as done here with *Psychotria*, excluding ranges that represent noncontiguous or at least widely disjunct distributions, as well as ranges that include a large number of areas, on the grounds that such ranges are implausible a priori.

### Alternative Models for Range Evolution

In the introductory section, we made the case that character models are ill-suited for biogeographic inference if single areas were considered equivalent to discrete character states, because a transition between states along a phylogenetic branch implies an unrealistic sequence of range evolution (dispersal followed by instantaneous local extinction in the source area). The DEC model also has limitations, such as its assumption that rates of geographic range evolution and speciation occur independently, and so it may not be appropriate for some kinds of biogeographic inference. In particular (and somewhat ironically), it may in fact be a poor fit to data on island biogeography as exemplified by the case of Hawaiian *Psychotria*; i.e., where terminal phylogenetic units represent populations, which by definition are restricted to single areas (islands). One might surmise that under such conditions, successful establishment of a new population by dispersal to another area actually represents a lineage divergence event, not merely range expansion, as in the DEC model. The rationale for this view lies in the assumption that allopatric lineage coalescence is rapid (effectively instantaneous with respect to the phylogenetic timescale) following dispersal to a new area.

Comparisons of empirical inferences about *Psychotria* using DEC models support this view. Under the preferred model M2, ancestral ranges are small and contiguous, and no local extinctions are inferred along branches, as reflected in the low estimated rate of 0.009 (Table 1). Five out of the six inferred dispersal events precede between-area cladogenesis events (Fig. 3c), a pattern consistent with dispersal being the proximate cause of lineage divergence. Taken together, the evidence is suggestive of a biogeographic history in which widespread

ancestors do not persist for any significant amount of time before diverging in allopatry.

A cladogenesis-by-dispersal model would differ from the DEC model in restricting inferences of ancestral ranges to single areas. It would differ from both DEC models and character models in constraining range evolution events (dispersal) to occur only at phylogenetic nodes, because in those cases, dispersal would be interpreted as the immediate cause of lineage divergence. Enumeration of scenarios involving the ancestral (source) area, retained by one daughter, and the destination area of the other dispersed daughter, would then become the ancestral states of interest. Along phylogenetic branches, “hidden” dispersal events yielding only extinct and unobserved descendants would need to be accounted for when estimating rates of dispersal, extinction, and cladogenesis. Fuller exploration of this model and DEC models await future papers. Clearly, model-based inference in historical biogeography is still in its infancy, with much theoretical work remaining.

### ACKNOWLEDGMENTS

We thank M. Nepokroeff for providing the phylogeny of *Psychotria*. Valuable feedback was gained from the working group on “Phytogeography of the Northern Hemisphere” sponsored by NESCent (NSF EF-0423641), S. Otto, and the Field Museum’s Evolution Discussion Group, especially N. Cordiero and A. Hipp. J. Clark, an anonymous reviewer, A. Baker, and J. Sullivan provided constructive comments on the original manuscript. RHR was partially supported by NSF grant DEB-0614108 and SAS was partially supported by NSF grant EF-0331654.

### REFERENCES

- Bremer, K. 1992. Ancestral areas: A cladistic reinterpretation of the center of origin concept. *Syst. Biol.* 41:436–445.
- Bremer, K. 1995. Ancestral areas: Optimization and probability. *Syst. Biol.* 44:255–259.
- Carson, H. L., and D. A. Clague. 1995. Geology and biogeography of the Hawaiian Islands. Pages 12–29 in *Hawaiian biogeography: Evolution on a Hot Spot Archipelago* (W. L. Wagner and V. A. Funk, eds.). Smithsonian Institution, Washington, DC.
- Edwards, A. W. F. 1992. *Likelihood*. Johns Hopkins University Press, Baltimore.
- Felsenstein, J. 1981. Evolutionary trees from DNA sequences: A maximum likelihood approach. *J. Mol. Evol.* 17:368–376.
- Funk, V. A., and W. L. Wagner. 1995. Biogeographic patterns in the Hawaiian Islands. Pages 379–419 in *Hawaiian biogeography: Evolution on a Hot Spot Archipelago* (W. L. Wagner and V. A. Funk, eds.). Smithsonian Institution, Washington, DC.
- Hausdorf, B. 1998. Weighted ancestral area analysis and a solution of the redundant distribution problem. *Syst. Biol.* 47:445–456.
- Huelsenbeck, J. P. and J. P. Bollback. 2001. Empirical and hierarchical Bayesian estimation of ancestral states. *Syst. Biol.* 50:351–366.
- Huelsenbeck, J., R. Nielsen, and J. Bollback. 2003. Stochastic mapping of morphological characters. *Syst. Biol.* 52:131–158.
- Maddison, W. P., P. E. Midford, and S. P. Otto. 2007. Estimating a binary character’s effect on speciation and extinction. *Syst. Biol.* 56:701–710.
- Nepokroeff, M., K. J. Sytsma, W. L. Wagner, and E. A. Zimmer. 2003. Reconstructing ancestral patterns of colonization and dispersal in the Hawaiian understory tree genus *Psychotria* (Rubiaceae): A comparison of parsimony and likelihood approaches. *Syst. Biol.* 52:820–838.
- Nielsen, R. 2002. Mapping mutations on phylogenies. *Syst. Biol.* 51:729–739.

- Pagel, M. 1994. Detecting correlated evolution on phylogenies: A general method for the comparative analysis of discrete characters. *Proc. R. Soc. Lond.* 255B:37–45.
- Pagel, M. 1999. The maximum likelihood approach to reconstructing ancestral character states of discrete characters on phylogenies. *Syst. Biol.* 48:612–622.
- Pagel, M., and A. Meade. 2006. Bayesian analysis of correlated evolution of discrete characters by reversible-jump Markov chain Monte Carlo. *Am. Nat.* 167:808–825.
- Pagel, M., A. Meade, and D. Barker. 2004. Bayesian estimation of ancestral character states on phylogenies. *Syst. Biol.* 53:673–684.
- Paradis, E. 2005. Statistical analysis of diversification with species traits. *Evolution* 59:1–12.
- Ree, R. H. 2005. Detecting the historical signature of key innovations using stochastic models of character evolution and cladogenesis. *Evolution* 59:257–265.
- Ree, R. H., B. R. Moore, C. O. Webb, and M. J. Donoghue. 2005. A likelihood framework for inferring the evolution of geographic range on phylogenetic trees. *Evolution* 59:2299–2311.
- Ronquist, F. 1994. Ancestral areas and parsimony. *Syst. Biol.* 43:267–274.
- Ronquist, F. 1997. Dispersal-vicariance analysis: a new approach to the quantification of historical biogeography. *Syst. Biol.* 45:195–203.
- Schluter, D., T. D. Price, A. O. Mooers, and D. Ludwig. 1997. Likelihood of ancestor states in adaptive radiation. *Evolution* 51:1699–1711.
- Slowinski, J. B. and C. B. Guyer. 1993. Testing whether certain traits have caused amplified diversification: An improved method based on a model of random speciation and extinction. *Am. Nat.* 142:1019–1024.

*First submitted 19 April 2007; reviews returned 5 June 2007;  
final acceptance 29 August 2007  
Associate Editor: Allan Baker*