

# Fast, Accurate and Simulation-Free Stochastic Mapping

Vladimir N. Minin<sup>1</sup> and Marc A. Suchard<sup>2,3</sup>

<sup>1</sup>Department of Statistics, University of Washington, Seattle, WA 98195-4322, USA

<sup>2</sup>Departments of Biomathematics and Human Genetics

David Geffen School of Medicine

<sup>3</sup>Department of Biostatistics, School of Public Health

University of California, Los Angeles, CA 90095-1766, USA

## Abstract

Mapping evolutionary trajectories of discrete traits onto phylogenies receives considerable attention in evolutionary biology. Given trait observations at the tips of a phylogenetic tree, researchers are often interested where on the tree the trait changes its state and whether some changes are preferential in certain parts of the tree. In a model-based phylogenetic framework, such questions translate into characterizing probabilistic properties of evolutionary trajectories. Moreover, if one adopts a Bayesian point of view, it is possible to incorporate uncertainty about model parameters, including phylogenies, into calculations of evolutionary trajectory properties. Current methods of assessing these properties rely on computationally expensive simulations. In this paper, we show that analytic approaches exist for two important and ubiquitous evolutionary trajectory properties. The first is the mean number of trait changes during evolutionary history, where changes can be divided into classes of interest (e.g. synonymous/nonsynonymous mutations). The mean evolutionary reward, accrued proportionally to the time a trait occupies each of its states, is the second property. We present an efficient, simulation-free algorithm for computing these two properties. Our exact algorithm requires only two tree traversals. Therefore, its computational efficiency is comparable to the familiar pruning algorithm for phylogenetic likelihood calculation. We provide two examples that illustrate practical importance of our method for Bayesian hypothesis testing. First, we employ posterior predictive diagnostics to detect correlation between two evolutionary traits. In the second example, we map synonymous and nonsynonymous mutations onto an HIV intrahost phylogeny and develop a formal test to compare fractions of nonsynonymous mutations on terminal and internal branches of an HIV intrahost phylogeny.

## Introduction and Background

Reconstructing evolutionary histories from present-day observations is a central problem in quantitative biology. Phylogenetic estimation is one example of such reconstruction. However, phylogenetic reconstruction alone does not provide a full picture of an evolutionary history, because evolutionary paths (mappings) describing trait states along the phylogenetic tree remain hidden. Although one is rarely interested in detailed reconstruction of such mappings, certain probabilistic properties of the paths are frequently used in evolutionary hypotheses testing (Nielsen, 2002; Huelsenbeck et al., 2003; Leschen and Buckley, 2007). For example, given a tree and a Markov model of amino acid evolution, one can compute the expected number of times a transition from a hydrophobic to a hydrophilic state occurs, conditional on the observed amino acid sequence alignment. Such expectations can inform researchers about model adequacy and provide insight into features of the evolutionary process overlooked by standard phylogenetic techniques (Dimmic et al., 2005).

Nielsen (2002) introduces stochastic mapping of trait states on trees and employs this new technique in a model-based evolutionary hypothesis testing context. The author starts with a discrete evolutionary trait  $X$  that attains  $m$  states. He further assumes that this trait evolves according to an evolutionary model described by a parameter vector  $\theta$ , where  $\theta$  consists of a tree  $\tau$  with  $n$  tips and branch lengths  $\mathbf{T} = (t_1, \dots, t_{B_n})$ , root distribution  $\pi = (\pi_1, \dots, \pi_m)$ , and a continuous-time Markov chain (CTMC) generator  $\mathbf{\Lambda} = \{\lambda_{ij}\}$  for  $i, j = 1, \dots, m$ . Let mapping  $\mathbf{M}_\theta = (\{X_{1t}\}, \dots, \{X_{B_n t}\})$  be a collection of CTMC trajectories along all branches of  $\tau$  and  $H(\mathbf{M}_\theta)$  be a real-valued summary of  $\mathbf{M}_\theta$ . Clearly, even when parameters  $\theta$  are fixed,  $h(\mathbf{M}_\theta)$  remains a random variable. Nielsen (2002) proposes to test evolutionary hypotheses using prior and posterior expectations  $E[H(\mathbf{M}_\theta)]$  and  $E[H(\mathbf{M}_\theta) | \mathbf{D}]$ , where  $\mathbf{D} = (D_1, \dots, D_n)$  are trait values observed at the  $n$  tips of

$\tau$ . Since these expectations are deterministic functions of  $\boldsymbol{\theta}$  and  $\mathbf{D}$ , they can be used as discrepancy measures for posterior predictive p-value calculations (Meng, 1994; Gelman et al., 1996).

A major advantage of Nielsen’s stochastic mapping framework is its ability to account for uncertainty in model parameters, including phylogenies. A major limitation of stochastic mapping is its current implementation that relies on time consuming simulations. In describing his method for calculating  $E[H(\mathbf{M}_{\boldsymbol{\theta}})]$  and  $E[H(\mathbf{M}_{\boldsymbol{\theta}}) | \mathbf{D}]$ , Nielsen (2002) writes “*In general, we can not evaluate sums in Equations 5 and 6 directly, because the set [of all possible mappings] is not of finite size.*” However, the infinite number of possible mappings does not prevent one from explicitly calculating  $E[H(\mathbf{M}_{\boldsymbol{\theta}})]$  for some choices of  $H$ . For example, if

$$H(\mathbf{M}_{\boldsymbol{\theta}}) = \begin{cases} 1 & \text{if } \mathbf{M}_{\boldsymbol{\theta}} \text{ is consistent with } \mathbf{D} \\ 0 & \text{if } \mathbf{M}_{\boldsymbol{\theta}} \text{ is inconsistent with } \mathbf{D}, \end{cases} \quad (1)$$

then  $E[H(\mathbf{M}_{\boldsymbol{\theta}})] = \Pr(\mathbf{D})$ , the familiar phylogenetic likelihood that can be evaluated without simulations (Felsenstein, 2004). Therefore, hope remains that other choices of  $H$  may also permit evaluation of  $E[H(\mathbf{M}_{\boldsymbol{\theta}})]$  and  $E[H(\mathbf{M}_{\boldsymbol{\theta}}) | \mathbf{D}]$  without simulations.

In this paper, we consider a class of additive mapping summaries of the form

$$H(\mathbf{M}_{\boldsymbol{\theta}}) = \sum_{b \in \Omega} h(\{X_{bt}\}), \quad (2)$$

where  $h(\{X_{bt}\})$  is a one-dimensional summary of the Markov chain path along branch  $b$  and  $\Omega$  is an arbitrary subset of all branches of  $\tau$ . Moreover, we restrict our attention to the two most popular choices of function  $h$ . Let  $\mathcal{L} \subset \{1, \dots, m\}^2$  be a set of ordered index pairs that label transitions of trait  $X$ . For each Markov path  $\{X_t\}$  and interval  $[0, t)$ , we count the number of labeled transitions in this interval and arrive at

$$h_1(\{X_t\}) = N_{\mathcal{L}} - \text{number of state transitions (labeled by set } \mathcal{L}, \quad (3)$$

where we omit dependence on  $\theta$  and  $t$  for brevity. Although our second choice of  $h$  is more abstract, it is motivated by Huelsenbeck et al. (2003), who use Nielsen’s stochastic mapping algorithm to calculate the mean dwelling time of a trait in a particular state. Let  $\mathbf{w} = (w_1, \dots, w_m)$  be a set of rewards assigned to each trait state. Trait  $X$  is “rewarded” the amount  $t \times w_i$  for spending time  $t$  in state  $i$ . We obtain the total reward of Markov path  $\{X_t\}$  by summing up all rewards that  $X$  accumulates during interval  $[0, t)$ ,

$$h_2(\{X_t\}) = R_{\mathbf{w}} - \text{evolutionary reward defined by vector } \mathbf{w}. \quad (4)$$

To obtain dwelling times of  $X$  in a predefined set of trait states, we set  $w_i = 1$  if  $i$  belongs to the set of interest and  $w_i = 0$  otherwise.

For these two choices of function  $h$ , we provide an algorithm for exact, simulation-free computation of  $E[H(\mathbf{M}_\theta)]$  and  $E[H(\mathbf{M}_\theta) | \mathbf{D}]$ . Similar to phylogenetic likelihood calculations of  $\Pr(\mathbf{D})$ , this algorithm relies on the eigen-decomposition of  $\mathbf{\Lambda}$  and requires traversing  $\tau$ . Despite the restricted form of these summaries, our results cover nearly all current applications of stochastic mapping. We conclude with two applications of stochastic mapping illustrating the capabilities of exact computation. In our first example, we examine co-evolution of two binary traits and demonstrate that a previously developed simulation-based test of independent evolution can be executed without simulations. We then turn to a large codon Markov state-space, on which simulation-based stochastic mapping generally experiences severe computational limitations. Using our exact computations, we study temporal patterns of synonymous and nonsynonymous mutations in intrahost HIV evolution.

## Local, One Branch Calculations

In this section, we provide mathematical details needed for calculating expectations of stochastic mapping summaries on one branch of a phylogenetic tree. We first motivate the need for such local computations by making further analogies between calculations of  $\Pr(\mathbf{D})$  and expectations of stochastic mapping summaries. The additive form of  $H$  reduces calculation of  $\mathbb{E}[H(\mathbf{M}_\theta)]$  and  $\mathbb{E}[H(\mathbf{M}_\theta) | \mathbf{D}]$  to computing branch-specific expectations  $\mathbb{E}[h(\{X_{bt}\})]$  and  $\mathbb{E}[h(\{X_{bt}\}) | \mathbf{D}]$ . Recall that according to most phylogenetic models, trait  $X$  evolves independently on each branch of  $\tau$ , conditional on trait states at all internal nodes of  $\tau$ . This conditional independence is the key behind the dynamic programming algorithm that allows for efficient calculation of  $\Pr(\mathbf{D})$  (Felsenstein, 1981). For this likelihood calculation algorithm, it suffices to compute finite-time transition probabilities  $\mathbf{P}(t) = \{p_{ij}(t)\}$ , where

$$p_{ij}(t) = \Pr(X_t = j | X_0 = i), \quad i, \quad (5)$$

for arbitrary branch length  $t$ . Similarly, to obtain  $\mathbb{E}[h(\{X_t\})]$  and  $\mathbb{E}[h(\{X_t\}) | \mathbf{D}]$ , we require means of computing local expectations  $\mathbf{E}(h, t) = \{e_{ij}(h, t)\}$ , where

$$e_{ij}(h, t) = \mathbb{E}[h(\{X_t\}) 1_{\{X_t=j\}} | X_0 = i], \quad (6)$$

and  $1_{\{.\}}$  is the indicator function. After illustrating how to compute  $\mathbf{E}(N_{\mathcal{L}}, t)$  and  $\mathbf{E}(R_{\mathbf{w}}, t)$  without resorting to simulations, we provide an algorithm that efficiently propagates local expectations  $\mathbf{E}(h, t)$  and finite-time transition probabilities  $\mathbf{P}(t)$  along  $\tau$  to arrive at  $\mathbb{E}[h(\{X_t\})]$  and  $\mathbb{E}[h(\{X_t\}) | \mathbf{D}]$ .

## Expected Number of Labeled Markov Transitions

Abstracting from phylogenetics, let  $N_{\mathcal{L}}(t)$  count the number of labeled transitions of a CTMC  $\{X_t\}$  during time interval  $[0, t)$ . It follows from the theory of Markov chain-induced counting processes that

$$\mathbf{E}(N_{\mathcal{L}}, t) = \int_0^t e^{\Lambda z} \mathbf{\Lambda}_{\mathcal{L}} e^{\Lambda(t-z)} dz, \quad (7)$$

where  $\mathbf{\Lambda}_{\mathcal{L}} = \{\lambda_{ij} \times 1_{\{(i,j) \in \mathcal{L}\}}\}$  (Ball and Milne, 2005). Since most evolutionary models are locally reversible, we can safely assume that  $\mathbf{\Lambda}$  is diagonalizable with eigen-decomposition  $\mathbf{\Lambda} = \mathbf{U} \times \text{diag}(d_1, \dots, d_m) \times \mathbf{U}^{-1}$ , where eigenvectors of  $\mathbf{\Lambda}$  form the columns of  $\mathbf{U}$ ,  $d_1, \dots, d_m$  are the real eigenvalues of  $\mathbf{\Lambda}$ , and  $\text{diag}(d_1, \dots, d_m)$  is a diagonal matrix with elements  $d_1, \dots, d_m$  on its main diagonal. Such analytic or numeric diagonalization procedure permits calculation of finite-time transition probabilities  $\mathbf{P}(t) = \mathbf{U} \times \text{diag}(e^{d_1 t}, \dots, e^{d_m t}) \times \mathbf{U}^{-1}$ , needed for likelihood calculations (Lange, 2004). Minin and Suchard (2008) show that one can use the same eigen-decomposition of  $\mathbf{\Lambda}$  to calculate local expectations

$$\mathbf{E}(N_{\mathcal{L}}, t) = \sum_{i=1}^m \sum_{j=1}^m \mathbf{S}_i \mathbf{\Lambda}_{\mathcal{L}} \mathbf{S}_j I_{ij}(t), \quad (8)$$

where  $\mathbf{S}_i = \mathbf{U} \mathbf{E}_i \mathbf{U}^{-1}$ ,  $\mathbf{E}_i$  is a matrix with zero entries everywhere except at the  $ii$ -th entry, which is one, and

$$I_{ij}(t) = \begin{cases} te^{d_i t} & \text{if } d_i = d_j, \\ \frac{e^{d_i t} - e^{d_j t}}{d_i - d_j} & \text{if } d_i \neq d_j. \end{cases} \quad (9)$$

## Expected Markov Rewards

For the reward process  $R_{\mathbf{w}}(t)$ , we define a matrix cumulative distribution function  $\mathbf{V}(x, t) = \{V_{ij}(x, t)\}$ , where

$$V_{ij}(x, t) = \Pr(R_{\mathbf{w}}(t) \leq x, X_t = j | X_0 = i). \quad (10)$$

Neuts (1995) demonstrates that local reward expectations can be expressed as

$$\mathbf{E}(R_{\mathbf{w}}, t) = -\frac{d}{ds} \mathbf{V}^*(s, t) \Big|_{s=0}, \quad (11)$$

where

$$\mathbf{V}^*(s, t) = \int_0^\infty e^{-sx} d\mathbf{V}(x, t) = e^{[\mathbf{\Lambda} - \text{diag}(w_1, \dots, w_m)]s} t \quad (12)$$

is the Laplace-Stieltjes transform of  $\mathbf{V}(x, t)$ . It is easy to see that the matrix exponential in equation (12) satisfies the following differential equation

$$\frac{d}{dt} \mathbf{V}^*(s, t) = \mathbf{V}^*(s, t) [\mathbf{\Lambda} - \text{diag}(w_1, \dots, w_m)] s. \quad (13)$$

Differentiating this matrix differential equation with respect to  $s$ , exchanging order of integration, and evaluating both sides of the resulting equation at  $s = 0$ , we arrive at the differential equation for local expectations

$$\frac{d}{dt} \mathbf{E}(R, t) = \mathbf{E}(R, t) \mathbf{\Lambda} + e^{\mathbf{\Lambda}t} \text{diag}(w_1, \dots, w_m), \quad (14)$$

where  $\mathbf{E}(R, 0)$  is the  $m \times m$  zero matrix. Multiplication of both sides of equation (14) by integrating factor  $e^{-\mathbf{\Lambda}t}$  from the right and integration with respect to  $t$  produces solution

$$\mathbf{E}(R_{\mathbf{w}}, t) = \int_0^t e^{\mathbf{\Lambda}z} \text{diag}(w_1, \dots, w_m) e^{\mathbf{\Lambda}(t-z)} dz. \quad (15)$$

Similarity between equations (7) and (15) invites calculation of the expected Markov rewards via spectral decomposition of  $\mathbf{\Lambda}$ ,

$$\mathbf{E}(R_{\mathbf{w}}, t) = \sum_{i=1}^m \sum_{j=1}^m \mathbf{S}_i \text{diag}(w_1, \dots, w_m) \mathbf{S}_j I_{ij}(t). \quad (16)$$

In summary, formulas (8) and (16) provide a recipe for exact calculations of local expectations for the number of labeled transitions and rewards.



# Assembling Pieces Together over a Tree

## Notation for Tree Traversal

Let us label the internal nodes of  $\tau$  with integers  $\{1, \dots, n-1\}$  starting from the root of the tree. Recall that we have already arbitrarily labeled the tips of  $\tau$  with integers  $\{1, \dots, n\}$ . Let  $\mathcal{I}$  be the set of internal branches and  $\mathcal{E}$  be the set of terminal branches of  $\tau$ . For each branch  $b \in \mathcal{I}$ , we denote the internal node labels of the parent and child of branch  $b$  by  $p(b)$  and  $c(b)$  respectively. We use the same notation for each terminal branch  $b$  except  $p(b)$  is an internal node index, while  $c(b)$  is a tip index. Let  $\mathbf{i} = (i_1, \dots, i_{n-1})$  denote the internal node trait states. Then, the complete likelihood of unobserved internal node states and the observed states at the tips of  $\tau$  is

$$\Pr(\mathbf{i}, \mathbf{D}) = \pi_{i_1} \prod_{b \in \mathcal{I}} p_{i_{p(b)} i_{c(b)}}(t_b) \prod_{b \in \mathcal{E}} p_{i_{p(b)} D_{c(b)}}(t_b). \quad (17)$$

We form the likelihood of the observed data by summing over all possible states of internal nodes,

$$\Pr(\mathbf{D}) = \sum_{i_1=1}^m \cdots \sum_{i_{n-1}=1}^m \pi_{i_1} \prod_{b \in \mathcal{I}} p_{i_{p(b)} i_{c(b)}}(t_b) \prod_{b \in \mathcal{E}} p_{i_{p(b)} D_{c(b)}}(t_b). \quad (18)$$

Clearly, when data on the tips are not observed, the prior distribution of internal nodes becomes

$$\Pr(\mathbf{i}) = \pi_{i_1} \prod_{b \in \mathcal{I}} p_{i_{p(b)} i_{c(b)}}(t_b). \quad (19)$$

## Posterior Expectations Of Mapping Summaries

Consider an arbitrary branch  $b^*$  connecting parent internal node  $p(b^*)$  to its child  $c(b^*)$ .

First, we introduce restricted moments

$$\mathbb{E}[h(\{X_{b^*t}\}) \mathbf{1}_{\mathbf{D}}] = \mathbb{E}[h(\{X_{b^*t}\}) | \mathbf{D}] \times \Pr(\mathbf{D}). \quad (20)$$

The expectation (20) integrates over all evolutionary mappings consistent with  $\mathbf{D}$  on the tips of  $\tau$ . Invoking the law of total expectation and the definition of conditional probability, we deduce

$$\begin{aligned}
\mathbb{E}[h(\{X_{b^*t}\})1_{\mathbf{D}}] &= \mathbb{E}[h(\{X_{b^*t}\}) | \mathbf{D}] \times \Pr(\mathbf{D}) = \sum_{\mathbf{i}} \mathbb{E}[h(\{X_{b^*t}\}) | \mathbf{i}, \mathbf{D}] \times \Pr(\mathbf{i}, \mathbf{D}) \\
&= \sum_{\mathbf{i}} \mathbb{E}[h(\{X_{b^*t}\}) | i_{p(b^*)}, i_{c(b^*)}] \pi_{i_1} \prod_{b \in \mathcal{I}} p_{i_{p(b)}i_{c(b)}}(t_b) \prod_{b \in \mathcal{E}} p_{i_{p(b)}D_{c(b)}}(t_b) \quad (21) \\
&= \sum_{\mathbf{i}} e_{i_{p(b^*)}i_{c(b^*)}}(h, t_{b^*}) \pi_{i_1} \prod_{b \in \mathcal{I} \setminus \{b^*\}} p_{i_{p(b)}i_{c(b)}}(t_b) \prod_{b \in \mathcal{E} \setminus \{b^*\}} p_{i_{p(b)}D_{c(b)}}(t_b).
\end{aligned}$$

The last expression in derivation (21) illustrates that in order to calculate the posterior restricted moment (20) along branch  $b^* \in \mathcal{I}$ , we merely need to replace finite-time transition probability  $p_{i_{p(b^*)}i_{c(b^*)}}(t_{b^*})$  with local expectation  $e_{i_{p(b^*)}i_{c(b^*)}}(h, t_{b^*})$  in the likelihood formula (18). Similarly, if  $b^* \in \mathcal{E}$ , we substitute  $e_{i_{p(b^*)}D_{c(b^*)}}(h, t_{b^*})$  for  $p_{i_{p(b^*)}D_{c(b^*)}}(t_{b^*})$  in (18). Given matrices  $\mathbf{P}(t_b)$  for  $b \neq b^*$  and  $\mathbf{E}(h, t_{b^*})$ , we can sum over internal node states using Felsenstein's pruning algorithm to arrive at the restricted mean  $\mathbb{E}[h(\{X_{b^*t}\})1_{\mathbf{D}}]$  and then divide this quantity by  $\Pr(\mathbf{D})$  to obtain  $\mathbb{E}[h(\{X_{b^*t}\}) | \mathbf{D}]$ .

This procedure is efficient for calculating the posterior expectations of mapping summaries for one branch of  $\tau$ . However, in practice, we need to calculate mapping expectations over many branches and consequently, execute the computationally intensive pruning algorithm many times. Schadt et al. (1998) encounter a similar problem during differentiation of the likelihood with respect to branch lengths. These authors formalize an algorithm that allows for computationally efficient, repeated replacement of one of the finite-time transition probabilities with an arbitrary function of the corresponding branch length in equation (18). This algorithm finds informal use since the 1980s in pedigree analysis (Cannings et al., 1980) and PAUP (personal communication, J. Huelsenbeck).

Let  $\mathbf{F}_u = (F_{u1}, \dots, F_{um})$  be a vector of forward, often called partial or fractional, likeli-

hoods at node  $u$ . Element  $F_{ui}$  is the probability of the observed data at only the tips that descend from node  $u$ , given that the state of  $u$  is  $i$ . If  $u$  is a tip, then we initialize partial likelihoods via equation  $F_{ui} = 1_{\{i=D_u\}}$ . In case of missing or ambiguous data,  $D_u$  denotes the subset of possible trait states, and forward likelihoods are set to  $F_{ui} = 1_{\{i \in D_u\}}$ . During the first, upward traversal of  $\tau$ , we compute forward likelihoods for each internal node  $u$  using the recursion

$$F_{ui} = \left[ \sum_{j=1}^m F_{c(b_1)j} p_{ij}(t_{b_1}) \right] \times \left[ \sum_{j=1}^m F_{c(b_2)j} p_{ij}(t_{b_2}) \right], \quad (22)$$

where  $b_1$  and  $b_2$  are indices of the branches descending from node  $u$  and  $c(b_1)$  and  $c(b_2)$  are the corresponding children of  $u$ . Schadt et al. (1998) suggest recording the quantities in square brackets in equation (22) together with  $\mathbf{F}_u$ . Naturally, these directional likelihoods  $\mathbf{S}_b = (S_{b_1}, \dots, S_{b_m})$  arise through the recursion

$$S_{bi} = \sum_{j=1}^m F_{c(b)j} p_{ij}(t_b). \quad (23)$$

Finally, we define backward likelihoods  $\mathbf{G}_u = (G_{u1}, \dots, G_{um})$ , where  $G_{ui}$  is the probability of observing state  $i$  at node  $u$  together with other tip states on the subtree of  $\tau$  obtained by removing all lineages downstream of node  $u$ . A second, downward traversal of  $\tau$  yields  $\mathbf{G}_u$  given the precomputed  $\mathbf{S}_b$ . We initialize the traversal by setting backward likelihoods at the root  $\mathbf{G}_1 = \boldsymbol{\pi}$ . Other  $\mathbf{G}_u$  follow from the recursion

$$G_{ui} = \sum_{j=1}^m G_{p(b)j} S_{bj} p_{ji}(t_b), \quad (24)$$

where  $b$  is the branch leading from node  $u$  to its parent  $p(b)$ .

For each branch  $b^*$ , we can sandwich  $p_{ij}(t_{b^*})$  among the forward, directional and backward, forward likelihoods and write

$$\Pr(\mathbf{D}) = \sum_{i=1}^m \sum_{j=1}^m G_{p(b^*)i} S_{b^*i} p_{ij}(t_{b^*}) F_{c(b^*)j}, \quad (25)$$

where  $b'$  is the second branch descending from the parent node  $p(b^*)$ . Therefore, with  $\mathbf{F}_u$ ,  $\mathbf{S}_b$  and  $\mathbf{G}_u$  precomputed for all of  $\tau$ , we can replace  $p_{ij}(t_{b^*})$  with any other quantity for any arbitrary branch  $b$  without repeatedly traversing  $\tau$ . In particular, the posterior restricted moment for branch  $b^*$  can be expressed as

$$\mathbb{E}[h(\{X_{b^*t}\})1_{\mathbf{D}}] = \sum_{i=1}^m \sum_{j=1}^m G_{p(b^*)_i} S_{b'i} e_{ij}(h, t_{b^*}) F_{c(b^*)_j}. \quad (26)$$

In Figure 1, we use an example tree to illustrate the correspondence between each quantity in sandwich formula (26) and the part of the tree involved in this quantity computation.

Remarkably, traversing  $\tau$  twice and caching  $\mathbf{F}_u$ ,  $\mathbf{S}_b$  and  $\mathbf{G}_u$  allows one to calculate posterior expectations of global additive mapping summaries  $\mathbb{E}[h(\{X_{b^*t}\}) | \mathbf{D}]$  repeatedly without any further traversals. We summarize all steps that lead to the computation of the global mean  $\mathbb{E}[H(\mathbf{M}_{\theta}) | \mathbf{D}]$  in Algorithm 1. Notice that  $\Pr(\mathbf{D})$ , needed to transition between conditional and restricted expectations in formula (20), is computed with virtually no additional cost in step 4 of the algorithm.

## Pulley Principle for Evolutionary Mappings

Suppose that we are interested in a mean mapping summary  $\mathbb{E}[H(\mathbf{M}_{\theta}) | \mathbf{D}]$ , obtained as a sum of local mapping summaries over *all branches* of the phylogenetic tree  $\tau$ . We would like to know whether quantity  $\mathbb{E}[H(\mathbf{M}_{\theta}) | \mathbf{D}]$  changes when we move the root of  $\tau$  to a different location.

Recall that reversibility of the Markov chain  $\{X_t\}$  makes  $\Pr(\mathbf{D})$  invariant to the root placement in  $\tau$  if the root distribution  $\boldsymbol{\pi}$  is the stationary distribution of  $\{X_t\}$  (Felsenstein, 1981). Felsenstein's pulley principle rests on the detailed balance condition

$$\pi_i p_{ij}(t) = \pi_j p_{ji}(t) \quad (27)$$

and Chapman-Kolmogorov relationship

$$p_{ij}(t_1 + t_2) = \sum_{k=1}^m p_{ik}(t_1)p_{kj}(t_2). \quad (28)$$

Applying both formulas (27) and (28) once in equation (18) allows one to move the root to any position along the two root branches without changing  $\Pr(\mathbf{D})$ . Therefore,  $\Pr(\mathbf{D})$  is invariant to moving the root to any position on any branch of  $\tau$  since we can repeatedly apply the detailed balance condition and the Chapman-Kolmogorov equation.

Invariance of  $\Pr(\mathbf{D})$  to root placement together with formula (20) suggests that root position invariance of conditional expectations  $E[H(\mathbf{M}_\theta) | \mathbf{D}]$  holds if and only if invariance of joint expectations  $E[H(\mathbf{M}_\theta)1_{\mathbf{D}}]$  is satisfied. Consider a 2-tip phylogeny with branches of length  $t_1$  and  $t_2$  leading to observed trait states  $D_1$  and  $D_2$  respectively. According to formulas (2) and (21), we may expect that

$$\begin{aligned} E[H(\mathbf{M}_\theta)1_{\mathbf{D}}] &= \sum_{k=1}^m \pi_k e_{kD_1}(h, t_1) p_{kD_2}(t_2) + \sum_{k=1}^m \pi_k e_{kD_2}(h, t_2) p_{kD_1}(t_1) \\ &= \sum_{k=1}^m \pi_{D_1} e_{D_1k}(h, t_1) p_{kD_2}(t_2) + \sum_{k=1}^m \pi_{D_2} e_{kD_2}(h, t_2) p_{D_1k}(t_1) \\ &= \pi_{D_1} e_{D_1D_2}(h, t_1 + t_2) \end{aligned} \quad (29)$$

depends only on the sum  $t_1 + t_2$ . Therefore, we can move the root anywhere on this phylogeny without altering expectations. It is easy to see that repeated application of derivation (29) readily allows for extension of the root invariance principle to  $n$ -tip phylogenies.

In derivation (29), we use identities

$$e_{ij}(h, t_1 + t_2) = \sum_{k=1}^m [e_{ik}(h, t_1) p_{kj}(t_2) + e_{kj}(h, t_2) p_{ik}(t_1)] \quad (30)$$

and

$$\pi_i e_{ij}(h, t) = \pi_j e_{ji}(h, t). \quad (31)$$

Equation (30) splits computation of the expected summary on interval  $[0, t_1 + t_2)$  into calculations bound to intervals  $[0, t_1)$  and  $[t_1, t_1 + t_2)$  with the help of the total expectation

law and Markov property. The derivation parallels the derivation of Chapman-Kolmogorov equation (28). Identity (31) requires more care as it *does not hold* for all choices of function  $h$ . Using equation (16), detailed balance condition (31) holds for  $h_2 = R_{\mathbf{w}}$ . However equation (7) suggests that we can guarantee the detailed balanced condition (31) for  $h_1 = N_{\mathcal{L}}$  only when  $(i, j) \in \mathcal{L}$  if and only if  $(j, i) \in \mathcal{L}$ .

## Prior Expectations of Mapping Summaries

In many applications of stochastic mapping, one wishes to compare the prior to posterior expectations of summaries (Nielsen, 2002). In this section, we derive formulas necessary for computing prior expectations. Similar to our derivation of posterior expectations, we begin by considering an arbitrary branch  $b^*$  and use the law of total expectation to arrive at

$$\mathbf{E}[h(\{X_{b^*t}\})] = \begin{cases} \sum_{\mathbf{i}} e_{i_{p(b^*)}i_{c(b^*)}}(h, t_{b^*}) \pi_{i_1} \prod_{b \in \mathcal{I} \setminus \{b^*\}} p_{i_{p(b)}i_{c(b)}}(t_b) & \text{if } b^* \in \mathcal{I} \\ \sum_{\mathbf{i}} e_{i_{p(b^*)}}(h, t_{b^*}) \pi_{i_1} \prod_{b \in \mathcal{I}} p_{i_{p(b)}i_{c(b)}}(t_b) & \text{if } b^* \in \mathcal{E}, \end{cases} \quad (32)$$

where  $e_i(h, t) = \sum_{j=1}^m e_{ij}(h, t)$  is the marginal local expectation of the mapping summary.

Identity  $\mathbf{P}(t)\mathbf{1} = \mathbf{1}$  allows us to eliminate summation over some internal node states in formula (32) and consider only those internal nodes that lie on the path connecting the root of  $\tau$  and  $c(b^*)$ . If  $\boldsymbol{\pi}$  is the stationary distribution of  $\{X_t\}$ , then formula (32) together with identities  $\boldsymbol{\pi}^T \mathbf{P}(t) = \boldsymbol{\pi}^T$  and  $\mathbf{P}(t)\mathbf{1} = \mathbf{1}$  simplifies prior local expectations even further,

$$\mathbf{E}[h(\{X_{b^*t}\})] = \sum_i \sum_j \pi_i e_{ij}(h, t_{b^*}) = \boldsymbol{\pi}^T \mathbf{E}(h, t_{b^*}) \mathbf{1} = \begin{cases} \boldsymbol{\pi}^T \boldsymbol{\Lambda}_{\mathcal{L}} \mathbf{1} t_{b^*} & \text{if } h = N_{\mathcal{L}} \\ \sum_{i=1}^m \pi_i w_i t_{b^*} & \text{if } h = R_{\mathbf{w}}. \end{cases} \quad (33)$$

The fact that prior local expectations at stationarity compute with virtually no additional burden has immediate practical implications. In the context of the posterior predictive model checking, researchers often need to simulate  $L$  independent and identically distributed (iid) realizations  $\mathbf{D}_1, \dots, \mathbf{D}_L$  of data at the tips of  $\tau$  and then calculate

the summary-based discrepancy measure  $\frac{1}{L} \sum_{l=1}^L \mathbb{E} [h(\{X_t\}) | \mathbf{D}_l]$  (Nielsen, 2002). Since we know the “true” model under simulation, we can approximate this discrepancy measure with the prior expectation

$$\frac{1}{L} \sum_{l=1}^L \mathbb{E} [h(\{X_t\}) | \mathbf{D}_l] \approx \sum_{\mathbf{D}^*} \mathbb{E} [h(\{X_t\}) | \mathbf{D}^*] \Pr(\mathbf{D}^*) = \mathbb{E} [\{X_t\}], \quad (34)$$

where  $\mathbf{D}^*$  ranges over all possible trait values at the tips of  $\tau$ . In molecular evolution applications,  $L$  is on the order of  $10^3 - 10^5$ , and hence approximation (34) should work well.

### Higher Moments and Variance of Mapping Summaries

So far we discuss calculations only for the first moments of additive mapping summaries. Higher moments are instrumental in some applications (Zheng, 2001; Nielsen, 2002). We first point out that local, one-branch calculations (8) and (16) extend easily to higher moments

$$e_{ij}^k(h, t) = \mathbb{E} \left[ h^k(\{X_t\}) 1_{\{X_t=j\}} | X_0 = i \right]. \quad (35)$$

Minin and Suchard (2008) explain how to perform such computations for counting processes. One-branch calculations of higher moments of reward processes are analogous.

To calculate higher moments of additive mappings summarizes over  $\tau$ , we need expectations of mixed product terms. For example, expectations of the form

$$\mathbb{E} [h(\{X_{b^*t}\})h(\{X_{b't}\})1_{\mathbf{D}}] \quad (36)$$

should be computed for some or all possible branch pairs  $(b^*, b')$  in order to obtain the second moment of the summary. Performing a derivation parallel to (21), we can show that calculating mixed product expectations (36) requires replacing finite-time transition probabilities  $p_{i_{p(b^*)}i_{c(b^*)}}(t_{b^*})$  and  $p_{i_{p(b')}i_{c(b')}}(t_{b'})$  with the first moments  $e_{i_{p(b^*)}i_{c(b^*)}}(h, t_{b^*})$  and

$e_{i_{p(b')} i_{c(b')}}(h, t_{b'})$  respectively. Unfortunately, caching partial likelihoods for repeated calculations becomes less practical in this case, because one needs to calculate partial likelihoods for each node pair in  $\tau$ . Clearly, such calculations result in an algorithm with running time growing as  $O(n^2)$ , where  $n$  is the number of tips.

Alternatively, we propose combining our analytic one-branch calculations with internal node state simulations to calculate the conditional variance of additive mapping summaries. Using the law of total variance, we decompose the posterior mapping variance as

$$\text{Var}[H(\mathbf{M}_\theta) | \mathbf{D}] = \text{E}\{\text{Var}[H(\mathbf{M}_\theta) | \mathbf{i}, \mathbf{D}]\} + \text{Var}\{\text{E}[H(\mathbf{M}_\theta) | \mathbf{i}, \mathbf{D}]\}. \quad (37)$$

Since  $X_t$  evolves independently on each branch conditional on internal node states  $\mathbf{i}$  and  $\mathbf{D}$ ,  $\text{Var}[H(\mathbf{M}_\theta) | \mathbf{i}, \mathbf{D}]$  and  $\text{E}[H(\mathbf{M}_\theta) | \mathbf{i}, \mathbf{D}]$  can be calculated using one-branch calculations. Using Monte Carlo integration, we are able to compute the expectation of the conditional variance and the variance of the conditional expectation by simulating internal node states from their posterior distribution  $\text{Pr}(\mathbf{i} | \mathbf{D})$ .

## Comparison with Simulation-Based Stochastic Mapping

We comment earlier that the Monte Carlo algorithm for stochastic mapping is an alternative and very popular way to compute expectations of mapping summaries. This algorithm consists of two major steps. In the first step, the tree is traversed once to compute  $\mathbf{F}_u$  for each node. In the second step, internal node states  $\mathbf{i}$  are simulated conditional on  $\mathbf{D}$  (Pagel, 1999). Then, conditional on  $\mathbf{i}$ , one simulates CTMC trajectories on each branch of the tree and computes summaries of interest. This second step is repeated  $N$  times producing a Monte Carlo sample of mapping summaries whose averages approximate the branch-specific expectations.

The running time of the algorithm depends on  $N$  and the computational efficiency of



generating CTMC trajectories. The number of samples  $N$ , required for accurate estimation, varies with  $\mathbf{A}$  and  $\mathbf{T}$ . Unfortunately, this aspect of stochastic mapping is largely ignored in the literature and by practitioners. Although more than one way exists to simulate trajectories conditional on starting and ending states (Rodrigue et al., 2008), rejection sampling is the most common (Nielsen, 2002). Assessing the efficiency of rejection sampling is complicated, because the efficiency depends not only on the choice of CTMC parameters, but also on the observed data patterns  $\mathbf{D}$ . We illustrate these difficulties using a CTMC on a state space of 64 codon triplets. We take two sites from an alignment of 129 HIV sequences that we discuss later in one of our examples. We call the first site slow evolving as it obtains only 3 different codons. The second, fast evolving site emits 9 different codons. We take one parameter slice of our Markov chain Monte Carlo (MCMC) sample and run rejection sampling to estimate the expected number of mutations for all branches of  $\tau$ . For each trial, we record the total number of rejections on all branches of  $\tau$  and the absolute errors of the estimated mean number of synonymous mutations summed over all branches. We summarize the results of this experiment in Table 1. We see a 400 $\times$  increase in the number of rejections required to simulate trajectories conditional on the fast site compared to equivalent simulations based on the slow site. The Monte Carlo error decreases as the number of Monte Carlo iterations increases, but not as fast as one would hope.

In summary, simulation-based stochastic mapping requires simulation of CTMC trajectories; this is not a trivial computational task. Assessing accuracy of methods is cumbersome and difficult to automate. Our Algorithm 1 replaces both simulation components from stochastic mapping calculations and therefore should be a preferred way of calculating expectations of mapping summaries. For the variance of mapping summaries, despite the added simulation component, our proposed procedure should still be more efficient than

current approaches, because we employ exact one-branch calculations to remove one of the two simulation layers.

## Examples

### Detecting Co-Evolution via Dwelling Times

In this section, we reformulate a previously developed simulation-based method for detection of correlated evolution (Huelsenbeck et al., 2003) in terms of a Markov reward process. We consider two primate evolutionary traits, estrus advertisement (EA) and multimale mating system (MS), analyzed by Pagel and Meade (2006). These authors first use *cytochrome-b* molecular sequences to estimate the posterior distribution of the phylogenetic relationship among 60 Old World monkeys and ape species. Using 500 MCMC samples from the posterior distribution of phylogenetic trees, Pagel and Meade run another MCMC chain, this time with a reversible jump component, that explores a number of CTMC models involving EA and MS traits and assess the models' posterior probabilities to learn about EA/MS co-evolution. The authors find support in favor of a hypothesis stating that EA presence correlates with MS presence. The trait data are shown in Figure 3 together with a phylogenetic tree, randomly chosen from the posterior sample. While not the case for Pagel and Meade (2006), reversible jump MCMC for model selection can be difficult to implement, especially as the number of trait states grows. Methods that simply require data fitting under the null model are warranted. Consequentially, we revisit this dataset and apply posterior predictive model diagnostics to test the hypothesis of independent evolution between EA and MS traits.

Our null model assumes that EA and MS evolve independently as two 2-state Markov

chains  $X_t^{(1)}, X_t^{(2)} \in \{0, 1\}$ , where 0 and 1 respectively stand for trait absence and presence.

Let the infinitesimal generators of the EA and MS CTMCs be

$$\mathbf{\Lambda}_{\text{EA}} = \begin{pmatrix} -\alpha_0 & \alpha_0 \\ \alpha_1 & -\alpha_1 \end{pmatrix} \quad \text{and} \quad \mathbf{\Lambda}_{\text{MS}} = \begin{pmatrix} -\beta_0 & \beta_0 \\ \beta_1 & -\beta_1 \end{pmatrix}. \quad (38)$$

We form a product Markov chain  $Y_t = (X_t^{(1)}, X_t^{(2)})$  on the state space  $\{(0, 0), (0, 1), (1, 0), (1, 1)\}$

that keeps track of presence/absence of the two traits simultaneously assuming that they evolve independently. The generator of the product chain is obtained via the Kronecker sum ( $\oplus$ )

$$\mathbf{\Phi} = \mathbf{\Lambda}_{\text{MS}} \oplus \mathbf{\Lambda}_{\text{EA}} = \begin{pmatrix} -(\alpha_0 + \beta_0) & \beta_0 & \alpha_0 & 0 \\ \beta_1 & -(\alpha_0 + \beta_1) & 0 & \alpha_0 \\ \alpha_1 & 0 & -(\alpha_1 + \beta_0) & \beta_0 \\ 0 & \alpha_1 & \beta_1 & -(\alpha_1 + \beta_1) \end{pmatrix}. \quad (39)$$

The Kronecker sum representation extends to general finite state-space Markov chains and to an arbitrary number of independently evolving traits (Neuts, 1995). Computationally, this representation is advantageous, because eigenvalues and eigenvectors of a potentially high-dimensional product chain generator derive analytically from eigenvalues and eigenvectors of low-dimensional individual generators (Laub, 2004).

To test the independent evolution model fit via posterior predictive diagnostics, we need a discrepancy measure (Meng, 1994). Following Huelsenbeck et al. (2003), we employ mean dwelling times to form a discrepancy measure. Let  $Z = \sum_{b=1}^B t_b$  be the tree length of  $\tau$ . We define the mean dwelling times  $Z_i^{(1)}$  and  $Z_i^{(2)}$  of traits  $X_t^{(1)}$  and  $X_t^{(2)}$  in state  $i$ , and the mean dwelling time  $Z_{ij}$  of the product chain  $Y_t$  in state  $(i, j)$  for  $i, j = 0, 1$ . More formally, we set

$$Z_i^{(1)} = \text{E} \left[ R_{\mathbf{w}_i} \mid \mathbf{D}^{(1)} \right], \quad Z_i^{(2)} = \text{E} \left[ R_{\mathbf{w}_i} \mid \mathbf{D}^{(2)} \right], \quad \text{and} \quad Z_{ij} = \text{E} \left[ R_{\mathbf{w}_{ij}} \mid \mathbf{D}^{(1)}, \mathbf{D}^{(2)} \right], \quad (40)$$

where  $\mathbf{w}_1 = (1, 0)$ ,  $\mathbf{w}_2 = (0, 1)$ ,  $\mathbf{w}_{00} = (1, 0, 0, 0)$ ,  $\mathbf{w}_{01} = (0, 1, 0, 0)$ ,  $\mathbf{w}_{10} = (0, 0, 1, 0)$ ,  $\mathbf{w}_{11} = (0, 0, 0, 1)$ , and  $\mathbf{D}^{(1)}$ ,  $\mathbf{D}^{(2)}$  are observations of the two traits on the tips of  $\tau$ .

Using the dwelling times, we define “expected”  $\varepsilon_{ij}$  and “observed”  $\eta_{ij}$  fractions of time the two traits spend in states  $i$  and  $j$

$$\varepsilon_{ij} = \frac{Z_i^{(1)}}{Z} \times \frac{Z_j^{(2)}}{Z} \text{ and } \eta_{ij} = \frac{Z_{ij}}{Z}. \quad (41)$$

We use quotation marks because both quantities are not observed. Under the null hypothesis of independence,  $\varepsilon_{ij} = \eta_{ij}$  for all  $i, j = 0, 1$ . To quantify the deviation from the null seen in the data, we introduce the discrepancy measure

$$\Delta(\mathbf{\Lambda}_{\text{EA}}, \mathbf{\Lambda}_{\text{MS}}, \mathbf{D}) = \sum_{i=0}^1 \sum_{j=0}^1 (\varepsilon_{ij} - \eta_{ij})^2. \quad (42)$$

This measure implicitly depends on  $\tau$  and branch lengths  $\mathbf{T}$ . We account for this dependence and phylogenetic uncertainty by averaging our results over a finite sample from the posterior distribution of  $\tau$  and  $\mathbf{T}$ , obtained from the molecular sequence data.

We use the software package BayesTraits to accomplish this averaging and to produce a MCMC sample from the posterior distribution of  $\mathbf{\Lambda}_{\text{EA}}$  and  $\mathbf{\Lambda}_{\text{MS}}$  assuming the null model of independent evolution of the two traits (Pagel et al., 2004). Each iteration sample from the output of BayesTraits consists of  $(\tau, \mathbf{T}, \alpha_0, \alpha_1, \beta_0, \beta_1)$  drawn from their posterior distribution. Given these model parameters, we generate a new dataset  $\mathbf{D}^{\text{rep}}$  and compute the observed  $\Delta(\mathbf{\Lambda}_{\text{EA}}, \mathbf{\Lambda}_{\text{MS}}, \mathbf{D})$  and predicted  $\Delta(\mathbf{\Lambda}_{\text{EA}}, \mathbf{\Lambda}_{\text{MS}}, \mathbf{D}^{\text{rep}})$  discrepancies for each iteration. We then compare their marginal distributions by plotting their corresponding histograms (Figure 2). In this figure, we also plot the observed against predicted discrepancies to display the correlation between these two random variables. The apparent disagreement between observed and predicted discrepancies is a manifestation of poor fit of the independent model of evolution. The observed discrepancy consistently exceeds the

predicted quantity. To illustrate the performance of predictive diagnostics when the independent model fits data well, we simulate trait data under this model along one of the 500 *a posteriori* supported phylogenetic trees. The second row of Figure 2 depicts the result of the posterior model diagnostics applied to the simulated data. In contrast to the observed primate data, the simulated data do not exhibit disagreement between the observed and predicted discrepancies.

Disagreement between the observed and predicted discrepancies can be quantified using a tail probability, called a posterior predictive p-value,

$$\text{ppp} = \Pr(\Delta(\mathbf{\Lambda}_{\text{EA}}, \mathbf{\Lambda}_{\text{MS}}, \mathbf{D}^{\text{rep}}) > \Delta(\mathbf{\Lambda}_{\text{EA}}, \mathbf{\Lambda}_{\text{MS}}, \mathbf{D}) \mid \mathbf{D}, H_0), \quad (43)$$

where the tail probability is taken over the posterior distribution of the independent model.

In practice, given  $N$  MCMC iterations, one estimates posterior predictive p-values via

$$\text{ppp} \approx \frac{1}{N} \sum_{g=1}^N \mathbf{1}_{\{\Delta(\mathbf{\Lambda}_{\text{EA}}^{(g)}, \mathbf{\Lambda}_{\text{MS}}^{(g)}, \mathbf{D}^{\text{rep},g}) > \Delta(\mathbf{\Lambda}_{\text{EA}}^{(g)}, \mathbf{\Lambda}_{\text{MS}}^{(g)}, \mathbf{D})\}}, \quad (44)$$

where  $\mathbf{\Lambda}_{\text{EA}}^{(g)}, \mathbf{\Lambda}_{\text{MS}}^{(g)}$  are parameter values, realized at iteration  $g$ , and  $\mathbf{D}^{\text{rep},g}$  is a data set, simulated using these parameter values. Following this recipe, we estimate ppp for the primate data and the artificial data. The discrepancy between the “observed” and “predicted” discrepancies is reflected in a low  $\text{ppp} = 0.0128$ . In contrast, the  $\text{ppp} = 0.3139$ , for the simulated data, supporting agreement between the “observed” and “predicted” distributions of  $\Delta$ .

## Mapping Synonymous and Non-Synonymous Mutations

In this section, we consider the important task of mapping synonymous and non-synonymous mutations onto branches of a phylogenetic tree. Our point of departure is a recent ambitious analysis of HIV intrahost evolution by Lemey et al. (2007), who use sequence data originally

reported by Shankarappa et al. (1999). The authors attempt to estimate branch-specific synonymous and nonsynonymous mutation rates and then project these measurements onto a time axis. This projection enables them to relate the time evolution of selection processes with clinical covariates. Lemey et al. (2007) find fitting codon models computationally prohibitive in this case. Instead, they first fit a DNA model to the intrahost HIV sequences, obtain a posterior sample of phylogenies with branch lengths, and then use these phylogenies to fit a codon model to the same DNA sequence alignment. Instead of fitting two different models to the data, we propose to use just a DNA model and exploit mapping summaries to predict synonymous and nonsynonymous mutation rates.

Suppose we observe a multiple DNA sequence alignment  $\mathbf{D} = (\mathbf{D}_1, \dots, \mathbf{D}_L)$  of a protein-coding region with  $L$  sites and that all  $C = L/3$  codons are aligned to each other such that the coding region starts at site 1 of  $\mathbf{D}$  (in other words, there is no frame-shift). We assume that sites corresponding to all first codon positions,  $\mathbf{D}_1, \mathbf{D}_4, \dots, \mathbf{D}_{L-2}$ , evolve according to a standard HKY model with generator  $\mathbf{\Lambda}_{\text{HKY}}(\kappa_1, \boldsymbol{\pi}_1)$  where  $\kappa_1$  is the transition-transversion ratio and  $\boldsymbol{\pi}_1$  is the stationary distribution, both just for the first codon position appropriately constrained. Similarly, we define CTMC generators  $\mathbf{\Lambda}_{\text{HKY}}(\kappa_2, \boldsymbol{\pi}_2)$  and  $\mathbf{\Lambda}_{\text{HKY}}(\kappa_3, \boldsymbol{\pi}_3)$  for the other two codon positions with independent parameters. Assuming that all  $L$  nucleotide sites in  $\mathbf{D}$  evolve independently together with the 3 codon position HKY models induces a product Markov chain model on the space of codons  $(AAA, AAG, \dots, TTT)$ , where codons are arranged in lexicographic order with respect to our nucleotide order  $A < G < C < T$ . The generator of this product CTMC is

$$\mathbf{\Lambda}_{\text{codon}} = \mathbf{\Lambda}_{\text{HKY}}(\kappa_1, \boldsymbol{\pi}_1) \oplus \mathbf{\Lambda}_{\text{HKY}}(\kappa_2, \boldsymbol{\pi}_2) \oplus \mathbf{\Lambda}_{\text{HKY}}(\kappa_3, \boldsymbol{\pi}_3). \quad (45)$$

With this Markov chain on the codon space, we define a labeling  $\mathcal{L}(s)$  that contains all possible pairs of codons that translate into the same amino acid. All other codon pairs

are collected into a labeling set  $\mathcal{L}(n)$ . Clearly, transitions between elements of  $\mathcal{L}(s)$  constitute synonymous mutations and nonsynonymous mutations are represented by transitions between elements of  $\mathcal{L}(n)$ . In this manner, counting processes map synonymous and nonsynonymous mutations onto specific branches of  $\tau$ . We consider HIV sequences from patient 1 of Shankarappa et al. (1999) and approximate the posterior distribution of our DNA model parameters  $\Pr(\mathbf{A}_{\text{codon}}, \tau, \mathbf{T} | \mathbf{D})$  using MCMC sampling implemented in the software package BEAST (Drummond and Rambaut, 2007). The serially sampled HIV sequences permit us to estimate the branch lengths  $\mathbf{T}$  in units of clock time, months in this case. For each saved MCMC sample, we compute branch-specific rates of synonymous and nonsynonymous mutations,

$$r_b(s) = \frac{\frac{1}{C} \sum_{c=1}^C \mathbb{E} [N_{\mathcal{L}(s)}(t_b) | \mathbf{D}_{c:c+2}]}{t_b} \quad \text{and} \quad r_b(n) = \frac{\frac{1}{C} \sum_{c=1}^C \mathbb{E} [N_{\mathcal{L}(n)}(t_b) | \mathbf{D}_{c:c+2}]}{t_b}, \quad (46)$$

where we denote data at codon site  $c$  by  $\mathbf{D}_{c:c+2}$ . We also record the fraction  $\frac{r_b(n)}{r_b(s)+r_b(n)}$  of nonsynonymous mutations. Similar to Lemey et al. (2007), we summarize these measurements by projecting them on the time axis. To this end, we form a finite time grid and produce a density profile of the synonymous and nonsynonymous rates, and of the nonsynonymous mutation fractions for each time interval between grid points (Figure 5). Both synonymous and nonsynonymous rate density profiles are consistently bimodal across time. Interestingly, the modes also stay appreciably constant. The density profile of the nonsynonymous mutation fractions is multimodal and fairly complex. There is a considerable number of branches that exhibit strong negative  $\left(\frac{r_b(n)}{r_b(s)+r_b(n)} \sim 0\right)$  and  $\left(\frac{r_b(n)}{r_b(s)+r_b(n)} \sim 1\right)$  positive selection. For the vast majority of branches, the nonsynonymous mutation fraction has first a modest upward trend through time and then descends to lower values, consistent with other patterns of evolutionary diversity reported by Shankarappa et al. (1999).

Intrigued by the multimodality observed in Figure 5, we investigate this issue further.

Lemey et al. (2007) consider several branch categories in their analysis, e.g. internal and terminal. We decide to test whether differences exist between selection forces acting on internal and terminal branches of  $\tau$ . We define the fractions of nonsynonymous mutations on internal and terminal branches as

$$\rho_{\mathcal{E}}(\mathbf{\Lambda}_{\text{codon}}, \tau, \mathbf{T}, \mathbf{D}) = \frac{\sum_{b \in \mathcal{E}} \sum_{c=1}^C \mathbb{E} [N_{\mathcal{L}(n)}(t_b) | \mathbf{D}_{c:c+2}]}{\sum_{b \in \mathcal{E}} \sum_{c=1}^C \{ \mathbb{E} [N_{\mathcal{L}(s)}(t_b) | \mathbf{D}_{c:c+2}] + \mathbb{E} [N_{\mathcal{L}(n)}(t_b) | \mathbf{D}_{c:c+2}] \}} \quad (47)$$

and

$$\rho_{\mathcal{I}}(\mathbf{\Lambda}_{\text{codon}}, \tau, \mathbf{T}, \mathbf{D}) = \frac{\sum_{b \in \mathcal{I}} \sum_{c=1}^C \mathbb{E} [N_{\mathcal{L}(n)}(t_b) | \mathbf{D}_{c:c+2}]}{\sum_{b \in \mathcal{I}} \sum_{c=1}^C \{ \mathbb{E} [N_{\mathcal{L}(s)}(t_b) | \mathbf{D}_{c:c+2}] + \mathbb{E} [N_{\mathcal{L}(n)}(t_b) | \mathbf{D}_{c:c+2}] \}}. \quad (48)$$

We plot their posterior histograms in Figure 4. These histograms do not overlap, suggesting different fractions of nonsynonymous mutations for internal and external branches. To test this hypothesis more formally, we form a discrepancy measure

$$\Delta(\mathbf{\Lambda}_{\text{codon}}, \tau, \mathbf{T}, \mathbf{D}) = \rho_{\mathcal{E}}(\mathbf{\Lambda}_{\text{codon}}, \tau, \mathbf{T}, \mathbf{D}) - \rho_{\mathcal{I}}(\mathbf{\Lambda}_{\text{codon}}, \tau, \mathbf{T}, \mathbf{D}). \quad (49)$$

As in our previous example, we compare the “observed” discrepancy  $\Delta(\mathbf{\Lambda}_{\text{codon}}, \tau, \mathbf{T}, \mathbf{D})$  with the “expected” discrepancy  $\Delta(\mathbf{\Lambda}_{\text{codon}}, \tau, \mathbf{T}, \mathbf{D}^{\text{rep}})$ , where  $\mathbf{D}^{\text{rep}}$  is a multiple sequence alignment simulated under the codon partitioning model with parameters  $\mathbf{\Lambda}_{\text{codon}, \tau}$  and  $\mathbf{T}$ . Evoking approximation (34) and recalling that our model assumes the same substitution rates for each branch of  $\tau$ , we deduce that

$$\Delta(\mathbf{\Lambda}_{\text{codon}}, \tau, \mathbf{T}, \mathbf{D}^{\text{rep}}) \approx 0 \quad (50)$$

for all parameter values  $\mathbf{\Lambda}_{\text{codon}}, \mathbf{T}$ , and  $\tau$  and replicated data  $\mathbf{D}^{\text{rep}}$ . Plugging in our new discrepancy measures into equation (44), we find that  $\text{ppp} < 0.001$ . Therefore, our posterior predictive test suggests that there is significant heterogeneity of mutation rates among the branches of  $\tau$ .



## Discussion

In this paper, we develop a computationally efficient framework for mapping evolutionary trajectories onto phylogenies. Although we aim to keep this mathematical framework fairly general, our main interest with evolutionary mappings lies in computing the mean number of labeled trait transitions and the mean evolutionary reward that depends linearly on the time a trait occupies each of its states. These two mapping summaries have been the most promising building blocks for constructing statistical tests. Incidentally, the transition counts and occupancy times also form the minimal sufficient statistics of partially observed CTMCs (Guttorp, 1995).

We build upon our earlier work involving single branch calculations for Markov-induced counting processes (Minin and Suchard, 2008). In our extension, we introduce single branch calculations for evolutionary reward processes and devise algorithms to extend single branch calculations to mapping expectations of counting and reward processes onto branches across an entire phylogeny. Our main result generalizes Felsenstein’s pruning algorithm that forms the work-horse of modern phylogenetic computation. The generalized pruning algorithm warrants two comments about its efficiency for performing simulation-free stochastic mapping. A traditionally slow component of phylogenetic inference is the eigen-decomposition of the infinitesimal rate matrix. Fortunately, this decomposition finds immediate re-use in our algorithm to calculate posterior expectations of mappings. Second, the algorithm requires only two traversals of the phylogenetic tree, and is therefore at most 2 times slower than the standard likelihood algorithm. In practice, we find our algorithm is about 1.5 slower than the likelihood calculation. We achieve this advantage because during the second traversal  $n$  terminal branches are not visited. Finally, the Felsenstein’s algorithm analogy yields a pulley principle for stochastic mapping and reduction in computation for prior expectations.

Our examples demonstrate how our novel algorithm facilitates phylogenetic exploratory analysis and hypothesis testing. First, we use simulation-free stochastic mapping of occupancy times to re-implement Huelsenbeck et al. (2003)'s posterior predictive test of independent evolution. In our second example, we attempt to recover synonymous and nonsynonymous mutation rates without resorting to codon models and instead use an independent codon partitioning model. We overcome this gross model misspecification with stochastic mapping, find intriguing multi-modality of synonymous and nonsynonymous rates, and use a posterior predictive model check to test differences in selection pressures between terminal and internal branches. We stress that our predictions are only as good as the model we use. For example, the terminal/internal branch differences may be due a general bad fit of our purposely misspecified model to the intrahost HIV data. However, we find this scenario unlikely in light of the recent demonstration of excellent performance of codon partitioning models during analyses of protein coding regions (Shapiro et al., 2006).

Our examples illustrate importance of hypothesis testing in statistical phylogenetics. In recent years, it has become clear that an evolutionary analysis almost never ends with tree estimation. Importantly, phylogenetic inference enables evolutionary biologists to tackle scientific hypotheses, appropriately accounting for ancestry-induced correlation in observed trait values (Huelsenbeck et al., 2000; Pagel and Lutzoni, 2002). Several authors demonstrate that mapping evolutionary histories onto inferred phylogenies provides a convenient and probabilistically grounded basis for designing statistically rigorous tests of evolutionary hypotheses (Nielsen, 2002; Huelsenbeck et al., 2003; Dimmic et al., 2005). Unfortunately, this important statistical technique has been hampered by the high computational cost of stochastic mapping. Our general mathematical framework and fast algorithms should secure a central place for stochastic mapping in the statistical toolbox of evolutionary biologists.

## Acknowledgments

We would like to thank Philippe Lemey for sharing his HIV codon alignments. We are also grateful to Andrew Meade for his help on using the BayesTraits software. MAS is supported by an Alfred P. Sloan Research Fellowship and a John Simon Guggenheim Memorial Fellowship.

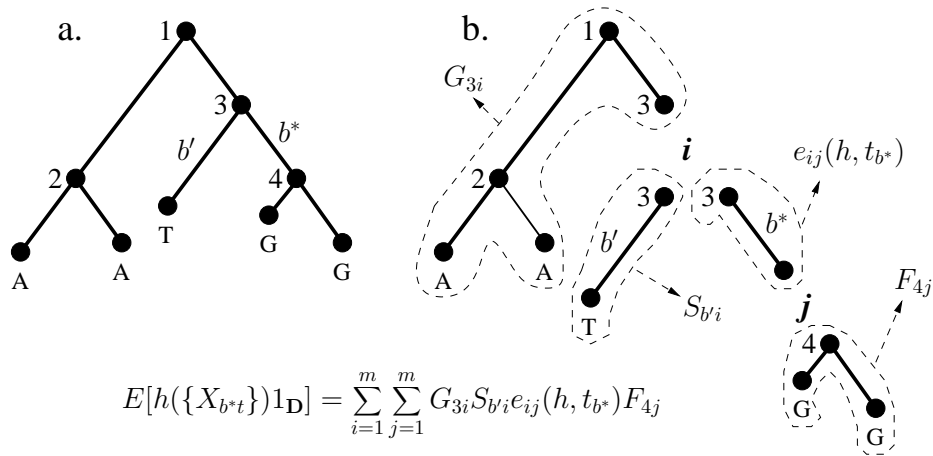
## References

- BALL, F. AND R. MILNE. 2005. Simple derivations of properties of counting processes associated with Markov renewal processes. *Journal of Applied Probability* 42:1031–1043.
- CANNINGS, C., E. THOMPSON, AND M. SKOLNICK. 1980. *Current Developments in Anthropological Genetics*, chapter Pedigree analysis of complex models, pages 251–298. Plenum Press, New York.
- DIMMIC, M., M. HUBISZ, C. BUSTAMANTE, AND R. NIELSEN. 2005. Detecting coevolving amino acid sites using Bayesian mutational mapping. *Bioinformatics* 21:i126–i135.
- DRUMMOND, A. AND A. RAMBAUT. 2007. BEAST: Bayesian evolutionary analysis by sampling trees. *BMC Evolutionary Biology* 7:214.
- FELSENSTEIN, J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach. *Journal of Molecular Evolution* 13:93–104.
- FELSENSTEIN, J. 2004. *Inferring Phylogenies*. Sinauer Associates, Inc., Sunderland, MA.
- GELMAN, A., X. MENG, AND H. STERN. 1996. Posterior predictive assessment of model fitness via realized discrepancies. *Statistica Sinica* 6:733–807.

- GUTTORP, P. 1995. *Stochastic Modeling of Scientific Data*. Chapman & Hall, Suffolk, Great Britain.
- HUELSENBECK, J., R. NIELSEN, AND J. BOLLBACK. 2003. Stochastic mapping of morphological characters. *Systematic Biology* 52:131–158.
- HUELSENBECK, J., B. RANNALA, AND J. MASLY. 2000. Accommodating phylogenetic uncertainty in evolutionary studies. *Science* 288:2349–2350.
- LANGE, K. 2004. *Applied Probability*. Springer-Verlag, New York.
- LAUB, A. 2004. *Matrix Analysis for Scientists and Engineers*. SIAM, Philadelphia, PA.
- LEMEY, P., S. K. POND, A. DRUMMOND, O. PYBUS, B. SHAPIRO, H. BARROSO, N. TAVEIRA, AND A. RAMBAUT. 2007. Synonymous substitution rates predict HIV disease progression as a result of underlying replication dynamics. *PLoS Computational Biology* 3:e29.
- LESCHEN, R. AND T. BUCKLEY. 2007. Multistate characters and diet shifts: Evolution of erotylidae (coleoptera). *Systematic Biology* 56:97–112.
- MENG, X. 1994. Posterior predictive p-values. *Annals of Statistics* 22:1142–1160.
- MININ, V. AND M. SUCHARD. 2008. Counting labeled transitions in continuous-time Markov models of evolution. *Journal of Mathematical Biology* 56:391–412.
- NEUTS, M. 1995. *Algorithmic Probability: a Collection of Problems*. Chapman and Hall, London.
- NIELSEN, R. 2002. Mapping mutations on phylogenies. *Systematic Biology* 51:729–739.

- PAGEL, M. 1999. The maximum likelihood approach to reconstructing ancestral character states of discrete characters on phylogenies. *Systematic Biology* 48:612–622.
- PAGEL, M. AND F. LUTZONI. 2002. *Biological Evolution and Statistical Physics*, chapter Accounting for phylogenetic uncertainty in comparative studies of evolution and adaptation, pages 148–161. Springer-Verlag, Berlin.
- PAGEL, M. AND A. MEADE. 2006. Bayesian analysis of correlated evolution of discrete characters by reversible–jump Markov chain Monte Carlo. *American Naturalist* 167:808–825.
- PAGEL, M., A. MEADE, AND D. BARKER. 2004. Bayesian estimation of ancestral character states on phylogenies. *Systematic Biology* 53:673–684.
- RODRIGUE, N., H. PHILIPPE, AND N. LARTILLOT. 2008. Uniformization for sampling realizations of Markov processes: applications to Bayesian implementations of codon substitution models. *Bioinformatics* 24:56–62.
- SCHADT, E., J. SINSHEIMER, AND K. LANGE. 1998. Computational advances in maximum likelihood methods for molecular phylogeny. *Genome Research* 8:222–233.
- SHANKARAPPA, R., J. MARGOLICK, S. GANGE, A. RODRIGO, D. UPCHURCH, H. FARZADEGAN, P. GUPTA, C. RINALDO, G. LEARN, X. HE, X. HUANG, AND J. MULLINS. 1999. Consistent viral evolutionary changes associated with the progression of human immunodeficiency virus type 1 infection. *Journal of Virology* 73:10489–10502.
- SHAPIRO, B., A. RAMBAUT, AND A. DRUMMOND. 2006. Choosing appropriate substitution models for the phylogenetic analysis of protein-coding sequences. *Molecular Biology and Evolution* 23:7–9.

ZHENG, Q. 2001. On the dispersion index of a Markovian molecular clock. *Mathematical Biosciences* 172:115–128.



$$E[h(\{X_{b^*t}\})1_{\mathcal{D}}] = \sum_{i=1}^m \sum_{j=1}^m G_{3i} S_{b'i} e_{ij}(h, t_{b^*}) F_{4j}$$

Figure 1: Sandwich formula illustration. In part (a), we plot an example phylogenetic tree in which we label internal nodes numerically and two branches  $b^*$  and  $b'$ . We break this tree at nodes 3 and 4 into the subtrees shown in part (b). Assuming that trait states are  $i$  and  $j$  at nodes 3 and 4 respectively, we mark each subtree by the corresponding quantity needed for calculating the posterior expectation of a mapping summary on branch  $b^*$ .

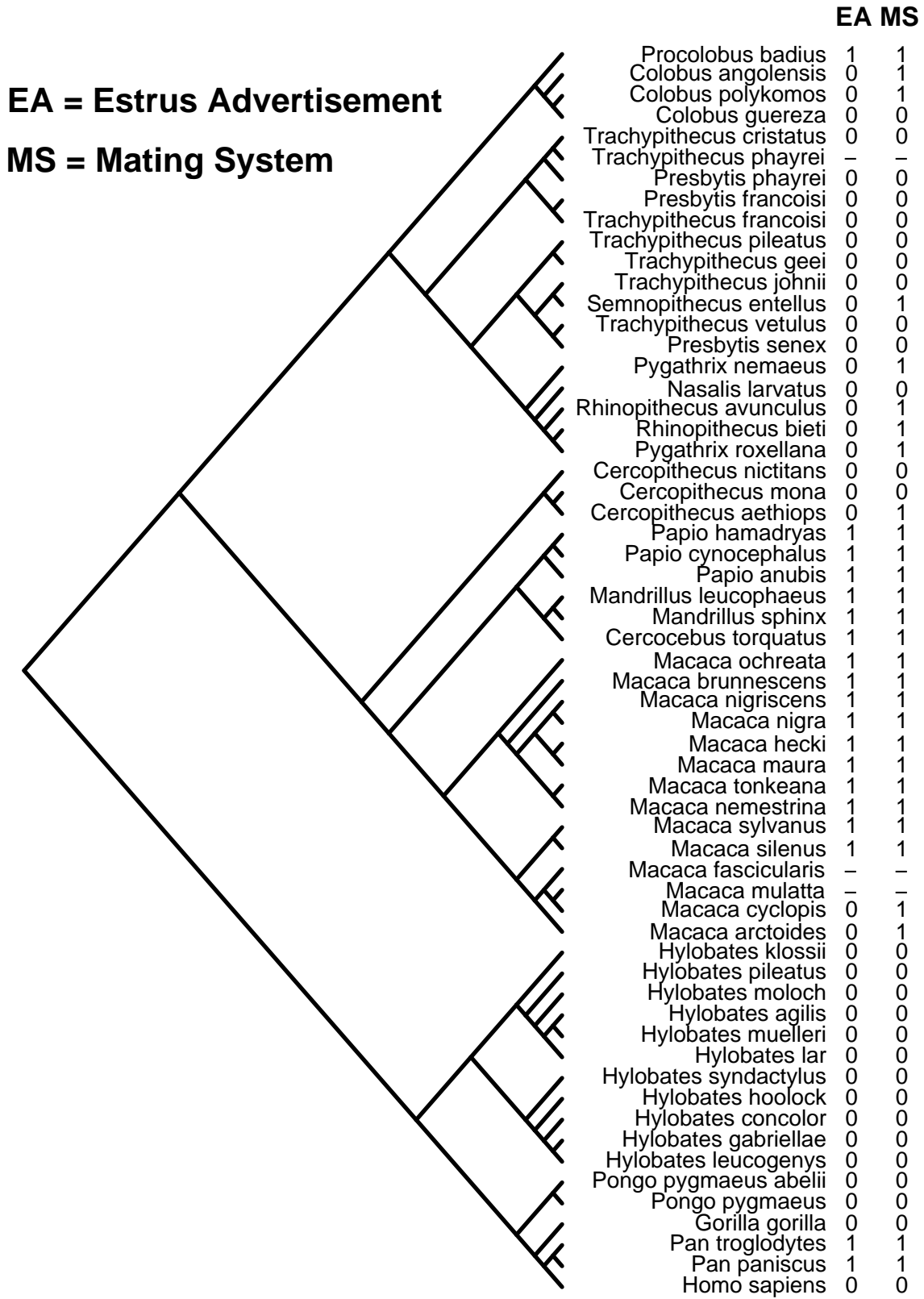


Figure 2: Primate trait data. We plot a phylogenetic tree, randomly chosen from the posterior sample, of 60 primate species. Branches of the tree are not drawn to scale. Taxa names and trait values (“0” - absence, “1” - presence, “-” - missing) for estrus advertisement (EA) and multi-male mating system (MS) are depicted at the tips of the tree.



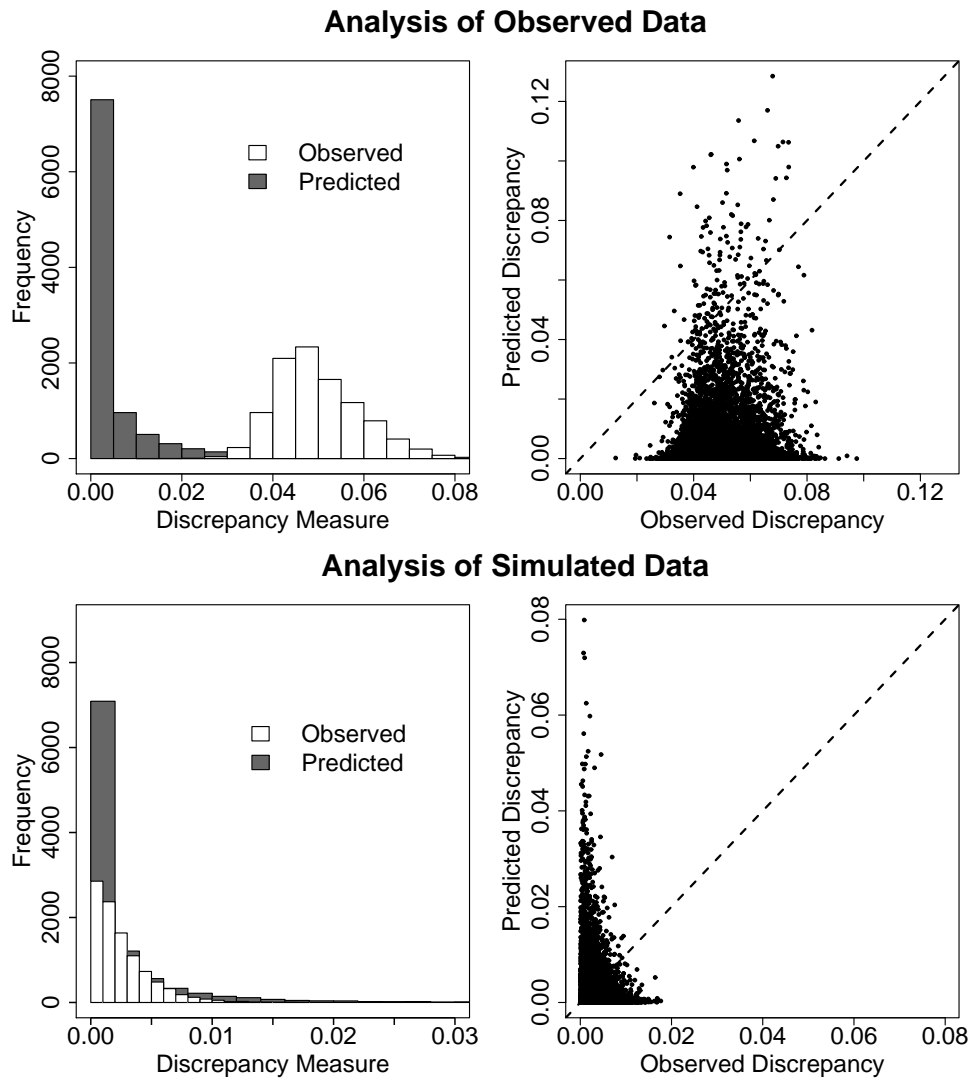


Figure 3: Testing co-evolution. The plots in the left column depict observed and predicted distributions of the discrepancy measure for the primate data (top) and simulated data (bottom). The right column shows the scatter plots of these distributions.

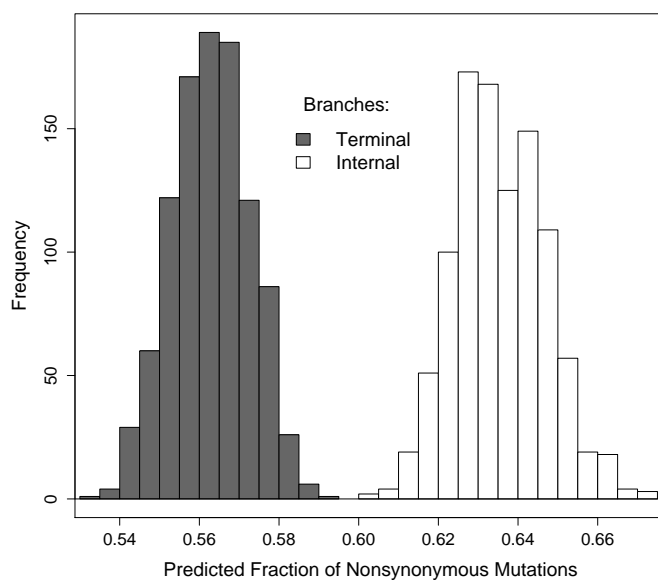


Figure 4: Bimodality of the fraction of nonsynonymous mutations. We plot the predicted fraction of nonsynonymous mutations computed for terminal (grey) and internal (white) branches.

## Intrahost HIV Evolution

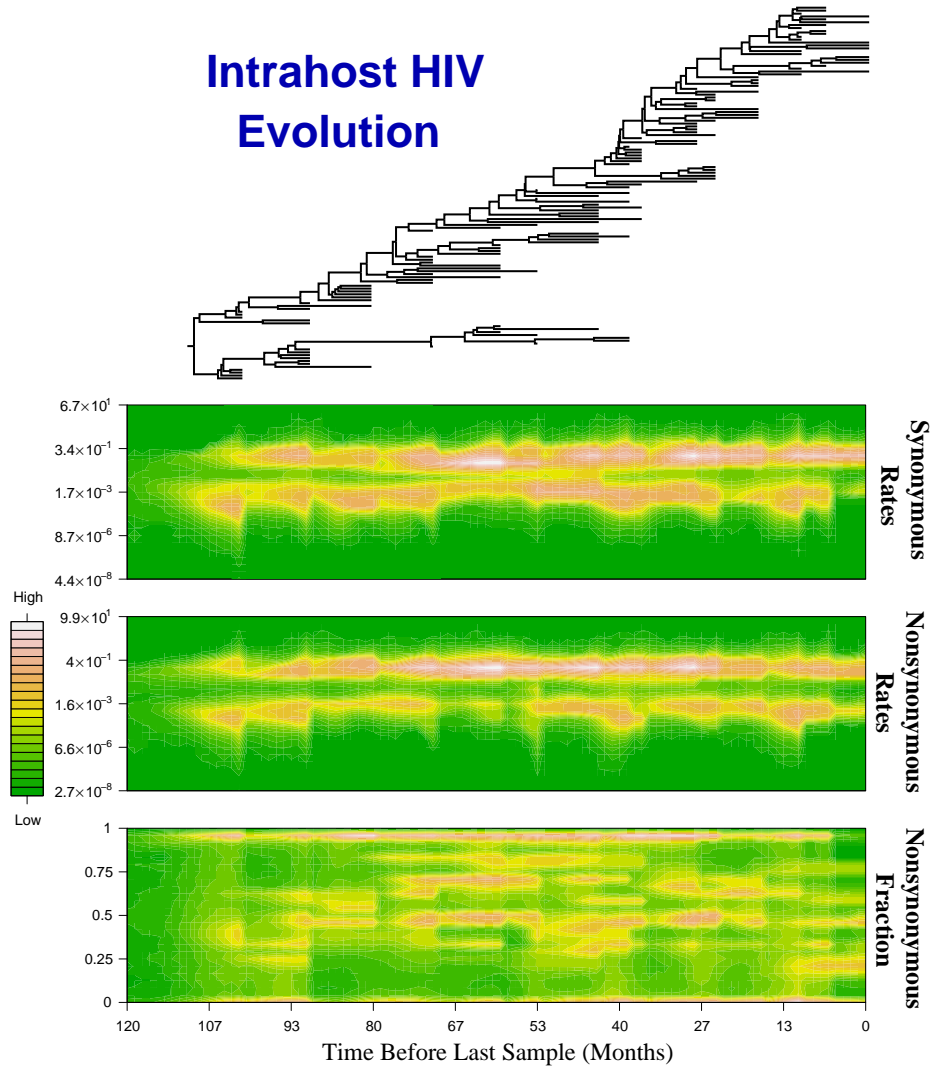


Figure 5: Time evolution of synonymous and nonsynonymous rates. At the top of the figure, we show a representative phylogeny of 129 intrahost HIV sequences. The three heat maps depict the marginal posterior densities of the synonymous and nonsynonymous rates, and the proportion of nonsynonymous mutations over time.

---

**Algorithm 1** Calculating Posterior Expectations  $\mathbf{E}[H(\mathbf{M}_\theta) | \mathbf{D}]$ 

---

- 1: Obtain an eigen-decomposition of the infinitesimal generator  $\mathbf{A}$
  - 2: Use this decomposition to compute  $\mathbf{P}(t_b)$  for each branch  $b$  of  $\tau$
  - 3: Employing the same eigen-decomposition, compute  $\mathbf{E}(h, t_{b^*})$  for each branch  $b^*$  in the set of interest  $\Omega$  using either equation (8) or equation (16)
  - 4: Traverse  $\tau$  once and use recursions (22) and (23) to calculate  $\mathbf{F}_u$  and  $\mathbf{S}_b$  for each node  $u$  and each branch  $b$ . Compute data likelihood  $\Pr(\mathbf{D})$  as the dot product of  $\mathbf{F}_{\text{root}}$  and root distribution  $\boldsymbol{\pi}$ .
  - 5: Traverse  $\tau$  the second time and calculate backward likelihoods  $\mathbf{G}_u$  for all nodes  $u$  via recursion equation (24)
  - 6: For each  $b^* \in \Omega$ , apply equation (26) to obtain  $\mathbf{E}[h(\{X_{b^*t}\})1_{\mathbf{D}}]$
  - 7: **return**  $\mathbf{E}[H(\mathbf{M}_\theta) | \mathbf{D}] = \frac{1}{\Pr(\mathbf{D})} \sum_{b^* \in \Omega} \mathbf{E}[h(\{X_{b^*t}\})1_{\mathbf{D}}]$
-

Table 1: Efficiency and accuracy of stochastic mapping. For each number of iterations, we report the median number of rejected CTMC trajectories over the entire phylogenetic tree per iteration and the sum of absolute errors (SAE) of simulation-based estimates of the mean number of synonymous mutations along branches of the phylogenetic tree.

Iterations	Slow Evolving Site		Fast Evolving Site	
	Rejections/Iteration	SAE	Rejections/Iteration	SAE
100	100	0.0598	38845	0.4624
500	105	0.0255	39247	0.3319
1000	102	0.0259	42075	0.2905
10000	106	0.0205	40805	0.2809