

Dated ancestral trees from binary trait data

Geoff K. Nicholls

Department of Statistics, Oxford, UK.

Russell D. Gray

Department of Psychology, Auckland University, Auckland, New Zealand.

Summary. Binary trait data record the presence or absence of distinguishing traits in individuals. We treat the problem of estimating ancestral trees with time depth from binary trait data. Birth times of traits visible in taxa are biased towards the leaves of the tree by the requirement that the trait survived to be observed. Unique traits, displayed at a single taxon, are commonly discarded. We model the evolution of binary traits as a birth-death process acting on the elements of sets of traits. The marginal prior distribution if the root time is uniform. We illustrate Bayesian inference for two binary-trait data sets which arise in historical linguistics. Model misspecification analysis is based in part on predictive distributions with external data, and related to distinctive features of the observation model.

Keywords: Phylogenetics, binary trait, dating methods, Bayesian inference, Markov chain Monte Carlo, glottochronology

1. Introduction

A great deal of progress has been made on the statistical analysis of DNA sequence data, and in particular for model-based estimation of genealogy. No equivalent statistical framework exists for trait-based cladistics. However, qualitative and quantitative trait data may be used to recover dated tree-like histories in situations where we have no genetic sequence data. Progress is possible when the traits are similar in type, so that some unifying assumption about their evolution is justified.

We give statistical methodology for tree-estimation from binary trait data. These data are made up of binary sequences, each sequence recording for one taxon the presence or absence of a list of traits. Pigeon wings and sparrow wings are *instances* of the trait “bird wings” displayed at the taxa “Pigeon” and “Sparrow”. In our model, two instances of a trait are necessarily *homologous*, that is, they descend from a common ancestor. Trait observation models have many missing data. Birth times of observed traits are unknown, and traits displayed at less than two taxa are commonly discarded. We model these missing data, integrate them out of the analysis analytically, and measure the random error using sample based Bayesian inference.

We begin in Section 2 with a model of the observation process. This model, described in Huson and Steel (2004) and Atkinson et al. (2005), is the natural stochastic process representing Dollo’s parsimony criterion, since each instance of a given trait descends from a single innovation. In Huson and Steel (2004) the traits are distinct genes which are

Address for correspondence: GK Nicholls, Department of Statistics, 1 South Parks Road, Oxford OX1 3TG, UK
E-mail: nicholls@stats.ox.ac.uk

present or absent in an individual, and trees are built using a maximum-likelihood pairwise-distance, and the neighbor-joining methods of Saitou and Nei (1987). Our own work has been motivated by a pair of data-sets, Dyen et al. (1997) and Ringe et al. (2002), recording trait presence and absence for Indo-European languages. Here traits are *cognate classes*, that is, homology classes of words of closely similar meaning. Thus English, Flemish and Danish share the trait “all/alle/al” whilst Spanish, Catalan and Italian lack that trait, but share “todo/tot/tutto”.

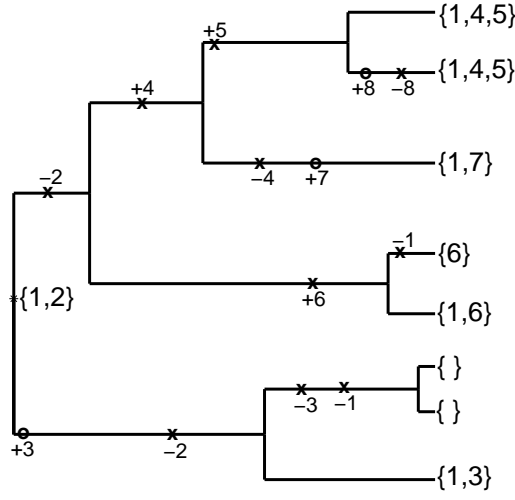
The survey given in Sankoff (1973) summarizes models of cognate trait data. Sankoff (1973) presents relatively realistic models which are complex and parameter-rich. Sankoff (1973) assumes inference will be based on pairwise distances between the binary-trait data-vectors of two languages. This mode of inference has been the norm for this data type. Thus Dyen et al. (1992) use classical hierarchical clustering of data-vectors based on pairwise distances between languages to establish a tree of languages. Gray and Atkinson (2003) use Ronquist and Huelsenbeck (2003) software and the Bayesian phylogenetic methods of Yang and Rannala (1997), to fit the finite-sites DNA sequence model of Felsenstein (1981). Pagel and Meade (2006) describe and fit a related, more realistic, model of cognate replacement within meaning category. These models allow traits identified in the data as homologous to arise by independent innovation. Warnow et al. (2006) propose a model in which each homology class has a unique birth event. However, there is to date no statistical inference for the model. Ringe et al. (2002), Erdem et al. (2005) and Nakhleh et al. (2005) reject dating, and avoid explicit modelling. They make a parsimony analysis without explicit measures of uncertainty. They allow some lateral transfer of traits, and thereby generalize to graphs which are not trees. They employ expert linguistic intervention in the inference, which becomes a well informed search through phylogenies. In light of the random and systematic error we measure below, we do not expect estimators of tree topology related to the mode (*ie* parsimony) to be adequate. However Ringe et al. (2002) add morphological traits. These may be more reliable data than cognate traits. Such traits can be analysed in the framework we set out.

Statistical contributions to the the dating of language branching events have been rejected by linguists. Dating efforts are criticized for their assumption of a constant, or at least homogeneous, rate of language change at all times and in all places, the so called “glottal clock”. We discuss this issue in detail in Section 9. Bergsland and Vogt (1962) found examples of extreme rates, but employed counter-examples biased by data-selection. Blust (2000) links rate heterogeneity to long-branch attraction. However, neither criticism considers even the random component of the error. In this respect we are repeating the comments of Sankoff (1973). Although we find evidence for model misspecification we nevertheless reproduce, to within random error, age estimates in analyses across near-independent data, and in reconstructions from synthetic data simulated under likely model-violation scenarios.

2. A model of binary trait evolution

In this section we specify an observation model for traits evolving on a fixed tree. Tree process models are discussed in Section 4.

We begin with notation for the tree. Let $g = (E, V, t)$ be a rooted binary tree with nodes $V = \{1, 2, \dots, 2L\}$, leaf nodes V_L , ancestral nodes V_A and edges $\langle i, j \rangle \in E$ where $i, j \in V$ and $i < j$. Node times $t = (t_1, t_2, \dots, t_{2L})$ are ordered $t_i \leq t_{i+1}$, so age increases towards the root. The root node label is $R = 2L - 1$ and there is an additional node $A = 2L$ with age



$L = 8$ leaves

$N = 4$ homologous trait classes in data
after thinning unique traits

$C = (c_1, c_2, c_3, c_4) = (1, 4, 5, 6)$

$M_1 = \{1, 2, 3, 5, 8\}$

$M_2 = \{1, 2\}$

$M_3 = \{1, 2\}$

$M_4 = \{4, 5\}$

$H(0, 1) = H(0, 2) = \{1, 4, 5\} = H_1 = H_2$

$H(0, 3) = \{1, 7\}, \quad H_3 = \{1\}$

$H(0, 4) = \{6\}, \quad H_4 = \{6\}$

$H(0, 5) = \{1, 6\}, \quad H_5 = \{1, 6\}$

$H(0, 6) = H(0, 7) = \{\}, \quad H_6 = H_7 = \{\}$

$H(0, 8) = \{1, 3\}, \quad H_8 = \{1\}$

$H(t_R, R) = \{1, 2\}$

Fig. 1. Observation process and notation. Trait birth (+c) and death (-c) events are marked. Birth events marked “o” generate traits absent from the data. Leaves are labelled $i = 1, 2, \dots, 8$ from top to bottom. NOUNIQUE is illustrated.

$t_A = \infty$ which is connected to the root via an edge $\langle R, A \rangle \in E$. Our convention is $A \notin V_A$, so $V = \{A\} \cup V_A \cup V_L$. Leaves may be staggered in time.

Next we give notation for the evolution of the sets of traits. The notation is illustrated in Fig. 1. Sets of trait labels evolve along the branches of g from the root towards the leaves, in the direction of decreasing age. Identify edge $\langle i, j \rangle$ by the node i at the base of that edge. For $\langle i, j \rangle \in E$ and $\tau \in [t_i, t_j]$ denote by (τ, i) a time point on a branch of g and by

$$[g] = \bigcup_{\langle i, j \rangle \in E} \bigcup_{\tau \in [t_i, t_j]} \{(\tau, i)\},$$

the set of all such points, including points on the edge $\langle R, A \rangle$ of infinite length. For each branch $\langle i, j \rangle \in E$ define a set-valued process $H(\tau, i) = \{h_1, h_2, \dots, h_{N(\tau, i)}\}$ of trait labels $h_a \in \mathbb{Z}$ for $a = 1, 2, \dots, N(\tau, i)$. The elements of $H(\tau, i)$ are realized by a simple reversible birth-death process which acts along each edge of the tree. Set elements (distinct traits) are born at constant rate λ and die at constant *per capita* rate μ . Set elements are distinguishable. At a branching event $(t_i, i) \in [g]$, set $H(t_i, i)$ is copied onto the top of the two branches $\langle j, i \rangle$ and $\langle k, i \rangle$ emerging from i , so that the evolution of $H(\tau, j)$ and $H(\tau, k)$ is conditional on $H(t_i, j) = H(t_i, k) = H(t_i, i)$.

The number of elements $N(t_R, R)$ in the root set is the number of traits born in $[t_R, \infty)$ which survive to time t_R . These surviving traits are generated at rate $\lambda(\tau, R) = \lambda \exp(-\mu(\tau - t_R))$, so their number is Poisson, mean λ/μ . The process may therefore be initialized by simulating $N(t_R, R) \sim \Pi(\lambda/\mu)$, and assigning $N(t_R, R)$ arbitrary trait labels $H(t_R, R) = \{1, 2, \dots, N(t_R, R)\}$ to the root set.

The data $D^{(L)} = (H_i, i \in V_L)$, with

$$H_i = H(t_i, i) \quad i \in V_L \quad (1)$$

are an ordered list of the sets of trait labels observed at the tree leaves. These data may be presented as a $L \times N$ binary matrix with data vectors (rows) corresponding to languages and variables (columns) corresponding to traits. Denote by $C = (\cup_{i \in V_L} H_i)$ the ordered set of all distinct trait labels appearing on the leaves, and let $C = (c_1, c_2, \dots, c_N)$. The matrix representation of the data is $D^{(LN)} = [D_{i,a}]_{i \in V_L}^{a=1..N}$ with $D_{i,a} = \mathbb{I}_{c_a \in H_i}$ (and \mathbb{I}_Z the indicator function for Z). We can represent the data as N sets of taxa labels also, with set M_a giving the leaves at which trait c_a appears. This representation is $D^{(N)} = (M_1, M_2, \dots, M_N)$, with $M_a = \{i : c_a \in H_i, i \in V_L\}$ for $a = 1..N$.

Traits displayed at just one taxon are often dropped from the data. It is argued that these unique traits do not inform tree topology. This is not the case in the model we have described, since unique traits are informative of time depth. Referring to the data analyzed in Section 7, Gray and Atkinson (2003) drop unique traits in their binary registration of the Dyen et al. (1997) data-set but retain them in their registration of the Ringe et al. (2002) data. The thinned data is $D^{(L)} = (H_i, i \in V_L)$, with

$$H_i = \left\{ c \in \mathbb{Z} : \sum_{i \in V_L} \mathbb{I}_{c \in H(t_i, i)} > 1 \right\} \quad i \in V_L \quad (2)$$

We call this observation model, which drops singleton traits, NOUNIQUE, in contrast to the model NOABSENT defined by Equation (1). We write D for generic NOABSENT or NOUNIQUE data. A realization of the NOUNIQUE observation process is shown in Fig. 1. Event +8 and the birth event +2 (which must occur on $\langle R, A \rangle$ so that 2 appears in the root-set $H(t_R, R)$) generate a trait which survives into no taxa at time $t = 0$. Events +3 and +7 generate a trait which survives into just a single taxon. For leaf i , the elements of the associated set $H(0, i)$ are instances of traits.

Felsenstein (1992) gives the likelihood for a Poisson process acting on a finite state space, along the branches of a tree, conditioned to show states other than the zero state at the leaves. The NOABSENT likelihood is similar. The model we have described resembles the Watterson (1975) infinite sites model, but here trait-death is in effect back-mutation, and similar too to the infinite alleles model of Kimura and Crow (1964), though the number of alleles is not random, whilst the number of traits is random.

3. Likelihood calculations

The likelihood for g, μ and λ is given in terms of the distribution of the point process of birth points for traits displayed in the data. Let $X = \{X_1, X_2, \dots, X_N\}$ be a random set of trait birth-points in $[g]$. The Poisson process generating X is obtained by thinning realizations of a constant rate process. Suppose a trait with label c is born at $z \in [g]$; let $O(z) = \sum_{i \in V_L} \mathbb{I}_{c \in H_i}$ give the number of taxa displaying trait c (after any thinning). If $\Pr\{O(z) > d | z, g, \mu\}$ is the probability for a trait, born at $z \in [g]$ to appear in the data at $d+1$ or more leaves, then the trait birth-rate at z in process X is $\lambda(z) = \lambda \Pr\{O(z) > d | z, g, \mu\}$, where $d = 0$ under the NOABSENT observation model and $d = 1$ under NOUNIQUE.

The distribution of X is defined on the space \mathcal{X} of all regular subsets $x \subset [g]$. For $f : [g] \rightarrow \mathfrak{R}$, define the integral $\int_{[g]} f(z) dz$ along tree branches by

$$\int_{[g]} f(z) dz = \sum_{\langle i, j \rangle \in E} \int_{t_i}^{t_j} f((\tau, i)) d\tau.$$

Now, suppose $X = x$ with $x = \{x_1, x_2, \dots, x_N\}$ and $x_a = (\tau_a, i_a)$ for $a = 1..N$, so that $x_a \in [g]$ identifies the point on the tree where trait c_a was born. Let $dx_a = dz$ at $x_a = z$. The density of the random set $X = x$, with respect to $dx = dx_1 dx_2 \dots dx_N$ on \mathcal{X} , is

$$f_X(x|g, \mu, \lambda) = \exp\left(-\int_{[g]} \lambda(z) dz\right) \prod_{a=1}^N \lambda(x_a).$$

We note in passing that $N \sim \Pi\left(\int_{[g]} \lambda(z) dz\right)$; the total number of distinct traits in the data is on average $\int_{[g]} \lambda(z) dz$.

Trait birth points are nuisance parameters, which we integrate out of the likelihood under the density f_X . Denote by $\Pr\{M_a = m_a | x_a, g, \mu, O(x_a) > d\}$ the probability for a trait, born at x_a , to be displayed at the leaves listed in set m_a and no others, conditional on being displayed in at least d leaves. The likelihood, $P(D|g, \mu, \lambda)$, is

$$\begin{aligned} P(D|g, \mu, \lambda) &= \int_{\mathcal{X}} P(D|x, g, \mu) f_X(x|g, \mu, \lambda) dx \\ &= \frac{e^{-\int_{[g]} \lambda(z) dz}}{N!} \prod_{a=1}^N \lambda \int_{[g]} \Pr\{M_a = m_a | x_a, g, \mu, O(x_a) > d\} \Pr\{O(x_a) > d | x_a, g, \mu\} dx_a. \end{aligned}$$

The outcome $\{M_a = m_a, O(x_a) > d\}$ is identical to the outcome $\{M_a = m_a\}$ for traits in the data, since those traits already satisfy the thinning condition $\text{card } m_a > d$ for each $a = 1, 2, \dots, N$. It follows that events in the data satisfy

$$\Pr\{M_a = m_a | x_a, g, \mu, O(x_a) > d\} = \frac{\Pr\{M_a = m_a | x_a, g, \mu\}}{\Pr\{O(x_a) > d | x_a, g, \mu\}},$$

and consequently the likelihood is

$$P(D|g, \mu, \lambda) = \frac{1}{N!} \exp\left(-\int_{[g]} \lambda(z) dz\right) \prod_{a=1}^N \lambda \int_{[g]} \Pr\{M_a = m_a | x_a, g, \mu\} dx_a. \quad (3)$$

We compute $\lambda \int_{[g]} \Pr\{O(z) > d | z, g, \mu\} dz$ and the factors $\lambda \int_{[g]} \Pr\{M_a = m_a | x_a, g, \mu\} dx_a$ using recursions related to the pruning recursion of Felsenstein (1981). We begin with $\lambda \int_{[g]} \Pr\{O(z) > d | z, g, \mu\} dz$. A birth at a generic point (τ, i) can be shifted to the child node, (t_i, i) ,

$$\Pr\{O(\tau, i) > d | (\tau, i), g, \mu\} = \Pr\{O(t_i, i) > d | (t_i, i), g, \mu\} \exp(-\mu(\tau - t_i)),$$

and the integral over $[g]$ reduced to a sum over contributions from edges:

$$\lambda \int_{[g]} \Pr\{O(z) > d | z, g, \mu\} dz = \frac{\lambda}{\mu} \sum_{\langle i, j \rangle \in E} \Pr\{O(t_i, i) > d | (t_i, i), g, \mu\} (1 - e^{-\mu(t_j - t_i)}). \quad (4)$$

We are interested in the cases $d = 0$ and $d = 1$. Let $u_i^{(d)} \equiv \Pr\{O(t_i, i) = d | (t_i, i), g, \mu\}$ so

$$\Pr\{O(t_i, i) > d | (t_i, i), g, \mu\} = \begin{cases} 1 - u_i^{(0)} & d = 0, \\ 1 - u_i^{(0)} - u_i^{(1)} & d = 1. \end{cases}$$

We give recursions for the $u_i^{(d)}$. Consider a pair of edges $\langle j, i \rangle, \langle k, i \rangle$ in E . Let $\delta_{i,j} = e^{-\mu(t_i - t_j)}$. The recursions

$$u_i^{(0)} = \left((1 - \delta_{i,j}) + \delta_{i,j} u_j^{(0)} \right) \left((1 - \delta_{i,k}) + \delta_{i,k} u_k^{(0)} \right) \quad (5)$$

$$u_i^{(1)} = \delta_{i,j} (1 - \delta_{i,k}) u_j^{(1)} + \delta_{i,k} (1 - \delta_{i,j}) u_k^{(1)} + \delta_{i,j} \delta_{i,k} (u_j^{(1)} u_k^{(0)} + u_j^{(0)} u_k^{(1)}) \quad (6)$$

are evaluated from $u_i^{(0)} = 0$ and $u_i^{(1)} = 1$ at leaves $i \in V_L$.

We need now to compute $\lambda \int_{[g]} \Pr\{M_a = m_a | x_a, g, \mu\} dx_a$ for generic trait patterns. Trait c_a is born into an edge ancestral to all the leaf nodes which display it, so the edges of g which contribute to the integral dx_a are those edges, E_a say, on the path to node A from the most recent common ancestor of the leaf nodes in m_a . Also, m_a is non-empty, so $\Pr\{M_a = m_a | (\tau, i), g, \mu\} = \Pr\{M_a = m_a | (t_i, i), g, \mu\} \exp(-\mu(\tau - t_i))$. We write the integral over $[g]$ in terms of a sum over contributions from edges:

$$\lambda \int_{[g]} \Pr\{M_a = m_a | x_a, g, \mu\} dx_a = \frac{\lambda}{\mu} \sum_{\langle i, j \rangle \in E_a} \Pr\{M_a = m_a | (t_i, i), g, \mu\} (1 - e^{-\mu(t_j - t_i)}). \quad (7)$$

Let $V_L^{(i)}$ be the set of leaf nodes in V descended from node i , including node i itself if i is a leaf node. For leaf sets m_a let $m_a^{(i)} = V_L^{(i)} \cap m_a$. Consider two edges $\langle j, i \rangle, \langle k, i \rangle$ in E . Events are independent down the two branches,

$$\Pr\{M_a^{(i)} = m_a^{(i)} | (t_i, i), g, \mu\} = \Pr\{M_a^{(j)} = m_a^{(j)} | (t_i, j), g, \mu\} \Pr\{M_a^{(k)} = m_a^{(k)} | (t_i, k), g, \mu\}$$

and moving from the top (t_i, j) to the bottom (t_j, j) of branch $\langle j, i \rangle$,

$$\Pr\{M_a^{(j)} = m_a^{(j)} | (t_i, j), g, \mu\} = \begin{cases} \delta_{i,j} \times \Pr\{M_a^{(j)} = m_a^{(j)} | (t_j, j), g, \mu\} & \text{if } m_a^{(j)} \neq \emptyset, \\ (1 - \delta_{i,j}) + \delta_{i,j} u_j^{(0)} & \text{if } m_a^{(j)} = \emptyset. \end{cases} \quad (8)$$

The recursion is evaluated from the leaves,

$$\Pr\{M_a^{(j)} = m_a^{(j)} | (t_j, j), g, \mu\} = \begin{cases} 1 & \text{if } j \text{ is a leaf and } m_a^{(j)} = \{j\}, \\ 0 & \text{if } j \text{ is a leaf and } m_a^{(j)} = \emptyset. \end{cases}$$

The recursion need not reach the leaves. It can be evaluated from nodes j satisfying $m_a^{(j)} = \emptyset$, using Equation (8) since $u_j^{(0)}$ is computed for the $\int_{[g]} \lambda(z) dz$ evaluation.

4. Prior models on trees

In this section we specify two families of probability distributions over trees, which we use to represent prior information concerning the phylogeny.

One tree prior we use is a branching process G_L with rate θ stopped at the instant of the L th branching event (counting the branching at the root). Denote by Γ the space of G_L -realizable trees and by dg the measure $\prod_{i \in V_A} dt_i$, with counting measure on topologies. The process G_L determines a density

$$f_G(g|\theta) \propto \theta^{L-1} \exp(-\theta|g|)$$

with respect to dg , where $|g|$ is the sum of all branch lengths, excluding the branch $\langle R, A \rangle$. The same functional form of the density is used when tree leaves are offset in time..

In the tree estimation problems we have encountered, the value of t_R , the root age, is of particular interest. Scientific hypotheses are expressed in statements of the form “ $t_R \in [t_{\min}, t_{\max}]$ ”. This motivates a prior which is non-informative with respect to such hypotheses. One prior which is strongly informative for t_R is the prior density $f(g|T) \propto \mathbb{I}_{t_R \leq T}$, the uniform distribution over all trees in Γ with root age smaller than T , a fixed upper limit. We find that, for trees with isochronous leaves at $t_1 = t_2 = \dots = t_L = 0$, the marginal distribution of t_R is t_R^{L-2} (for each $g \in \Gamma$ the topology-constrained volume integral $\int dt_{L+1} \dots dt_{2L-2} \propto t_R^{L-2}$). This prior represents a state of belief in which $\Pr\{t_R \in [T/2, T]\}$ is about 2^{L-1} times greater than $\Pr\{t_R \in [0, T/2]\}$. The marginal density of t_R in the prior

$$f_R(g|T) \propto t_R^{2-L} \mathbb{I}_{t_R \leq T}$$

is uniform in $[0, T]$.

Where calibration constraints are imposed, the prior $f_R(g|T)$ must be modified. Date calibration data is prior information about the timing of events in g , typically knowledge of tree topology and node times for the portion of the tree adjacent to the leaves. In the examples in Section 7, certain complete subtrees, called clades, are imposed, along with upper and lower bounds on the age of the most recent common ancestor of the leaf nodes in these clades. These constraints are represented as black bars in Fig. 5 and Fig. 8. For each such clade a list of leaf nodes and a clade-root age range is specified. An admissible tree displays the leaf nodes in the list as a clade and their most recent common ancestor falls inside the age range specified for the root of that clade.

Tree nodes in clades with clade root times bounded above by calibration constraints do not contribute a factor t_R to the tree-topology constrained volume integral $\int \prod_{i \in V_A \setminus \{R\}} dt_i$. The density f_R must be further modified to take into account non-isochronous leaf dates. The exact result is beyond us. However, if S is a list of free nodes, *ie* nodes $i \in V_A \setminus \{R\}$ which do not belong to root-bounded clades, and for node $i \in S$, s_i is the minimum time-value node i can achieve in an admissible tree, we found that

$$f_R(g|T) \propto \mathbb{I}_{t_R < T} \prod_{i \in S} (t_R - s_i)^{-1}$$

gives a reasonably flat marginal distribution for t_R when T is much greater than the maximum upper bound on clade root times. In the examples following Section 7, we summarize posterior distributions computed under tree priors $f = f_G$ and $f = f_R$ with clade calibration constraints.

We encounter data in which leaf node times are themselves subject to uncertainty. Calibration data on leaf node times allow leaf times to vary in a range, so that for each $i \in V_L$, $t_i \in [t_i^-, t_i^+]$. The leaf times $t_i, i \in V_L$ become missing data. In Section 7, the allowed range for leaf times is small compared to the time over which traits evolve. We take a prior uniform in $[t_i^-, t_i^+]$ for $t_i, i \in V_L$.

5. Posterior distributions

Our final expression for the likelihood is obtained by substituting Equation (4) and Equation (7) into Equation (3), and evaluating these terms using Equation (5) and Equation (8)

respectively. Multiplying that likelihood by the tree-prior f_G given in Section 4 we obtain the posterior distribution

$$\begin{aligned}
p(g, \mu, \lambda, \theta | D) d\theta d\lambda d\mu dg &\propto \exp\left(-\frac{\lambda}{\mu} \sum_{\langle i, j \rangle \in E} \Pr\{O(t_i, i) > d|(t_i, i), g, \mu\} (1 - e^{-\mu(t_j - t_i)})\right) \\
&\times \prod_{a=1}^N \sum_{\langle i, j \rangle \in E_a} \Pr\{M_a = m_a | (t_i, i), g, \mu\} (1 - e^{-\mu(t_j - t_i)}) \\
&\times \left(\frac{\lambda}{\mu}\right)^N \theta^{L-1} e^{-\theta|g|} d\theta d\lambda d\mu \prod_{i \in V_A} dt_i. \tag{9}
\end{aligned}$$

Equation (9) holds for tree prior f_G . Under tree prior $f_R(g|T)$ we drop parameter θ from the posterior and replace $f_G \propto \theta^{L-1} e^{-\theta|g|}$ with $f_R \propto t_R^{2-L} \mathbb{1}_{t_R < T}$.

Time scale is undetermined under scale invariant priors $p(\mu, \lambda, \theta) = (\mu\lambda\theta)^{-1}$. For $\rho > 0$, the transformation $t_1, \dots, t_R, \mu, \lambda \rightarrow t_1/\rho, \dots, t_R/\rho, \mu\rho, \lambda\rho, \theta\rho$ leaves $p(g, \mu, \lambda, \theta | D) dg d\lambda d\mu d\theta$ invariant, so it cannot be a proper distribution. The problem remains (for $\rho > t_R/T$) under tree prior f_R . Date calibration data described in Section 4 restricts the space of tree states Γ , and breaks the time-rescaling invariance. The posterior becomes proper.

The special case of two taxa (so, $0 = t_1 \leq t_2 \leq t_3$ with $t_R = t_3$ the tree height) and $d = 0$ is of interest for checking and debugging. The data, $D^{(L)} = (H_1, H_2)$, are two lists of trait instances, including traits present at just one leaf. Huson and Steel (2004) compute the MLE, $|g|^*$, for the total tree length $|g| = 2t_R - t_2 - t_1$ in a two leaf tree directly, using the reversibility of the birth-death process of traits between the two leaves, and conditioned on λ/μ known. In this way they motivate a new measure of the distance between two binary sequences as the pairwise maximum likelihood distance between the two sequences, estimated under this model. Let $n_1 = \text{card } H_1 \setminus H_2$, $n_2 = \text{card } H_2 \setminus H_1$ and $n_{12} = \text{card } (H_1 \cap H_2)$, so that $N = n_{12} + n_1 + n_2$. The data D amounts to n_1, n_2, n_{12} in the two leaf case. The likelihood for the two leaf case, computed from Equation (3), using Equation (4) and Equation (7), is

$$P(n_1, n_2, n_{12} | |g|, \lambda, \mu) \propto \left(\frac{\lambda}{\mu}\right)^N \exp\left(-\frac{\lambda}{\mu} [2 - e^{-\mu|g|}]\right) (1 - e^{-\mu|g|})^{n_1 + n_2} e^{-\mu|g|n_{12}}.$$

The MLE we obtain for $|g|$ given λ/μ is not enlightening but agrees with that given in Huson and Steel (2004). Maximizing $P(n_1, n_2, n_{12} | |g|, \lambda, \mu)$ over λ and $|g|$ simultaneously we obtain the simple result

$$|g|^* = \frac{1}{\mu} \log\left(1 + \frac{n_1 + n_2}{2n_{12}}\right) \tag{10}$$

subject to the constraint $|g|^* \geq t_2 - t_1$. This result expresses an exponential decline in percentage shared traits with increasing total age separation between the two taxa. It has been widely imposed for pairwise analyses of trait data. Swadesh (1952) fits a relation of this kind to lexical trait data.

The posterior distribution for $|g|$ given μ , which is available in closed form for the two leaf tree, is useful for debugging MCMC code. Taking priors $p(\lambda, \theta) = (\lambda\theta)^{-1}$ in Equation (9)

and integrating out λ and θ , we obtain,

$$p(|g| | \mu, n_1, n_2, n_{12}) \propto \frac{1}{\mu|g|} \left[\frac{e^{-\mu|g|}}{2 - e^{-\mu|g|}} \right]^{n_{12}} \left[\frac{1 - e^{-\mu|g|}}{2 - e^{-\mu|g|}} \right]^{n_1 + n_2}. \quad (11)$$

We see that μ and $|g|$ appear in the combination $\mu|g|$. When we consider large trees, and estimate μ , calibration constraints fixing clades in g separate this pair of variables.

6. Markov chain Monte Carlo

We work exclusively with the marginal posterior density $p(g, \mu | D)$. When the prior for λ is λ^{-1} , this variable is Gamma distributed in the posterior, and may be integrated. The same observation applied to θ , when we use the f_G prior. Sampling the posterior distribution $p(g, \mu | D)$ via Metropolis-Hastings Markov chain Monte Carlo is fairly straightforward, once efficient schemes for evaluating and updating the recursions, Equations (5), (6) and (8) have been implemented.

We use the tree operations described in Drummond et al. (2002). These include updates which alter the tree topology, updates which vary node times, updates which vary parameters such as μ , and updates which make some combination of these changes. In a specimen update we generate candidates for Metropolis-Hastings updates by simulating $\rho \sim U(1/2, 2)$ and setting $t' = \rho t$ and $\mu' = \mu/\rho$, since this is expected to be a ridge direction of the loglikelihood. In the acceptance probability for this update, the probability density to generate the reverse update, with $\rho' = 1/\rho$, is equal to the probability density to generate the forward update, and a Jacobian term $|\partial(g', \mu', \rho')/\partial(g, \mu, \rho)| = \rho^{L-2}$ appears in the Hastings ratio.

Our MCMC convergence analysis, based on monitoring the asymptotic behavior of the autocorrelation for μ , t_R , $|g|$ and the log-likelihood, follows Geyer (1992). The error bars in Figures 3, 11 and 14 were estimated in the following way. If $\tau^{(f)}$ is the integrated autocorrelation time, estimated using Geyer's monotone sequence estimator, then, in a run of length J , generating samples $\{X_i\}_{i=1}^J$, $\hat{\sigma}(\bar{f})^2 = \hat{\sigma}(f)^2 \hat{\tau}^{(f)}/J$ is an estimate of the variance of $\bar{f} = J^{-1} \sum_i f(X_i)$. Marginal posterior distributions for individual ages and rates were close to normal in shape.

We made a number of checks on our implementation. We check the likelihood sums to one over data. We check that the marginal prior distribution of t_R under f_R with isochronous leaves is uniform. We recover the posterior distribution in Equation (11) in the two leaf case. We fix a data set and vary the proportions in which update types are used. We check that statistics computed under the posterior do not vary, to within estimated errors. We recover the parameters of synthetic data, and the posterior distribution concentrates on the correct parameter values as the number, N , of traits displayed in the data increases.

7. Data

In the Dyen et al. (1997) and Ringe et al. (2002) data, a trait is a homology class of words. The setup is illustrated in Table 1. A set of K meaning categories are chosen and, for each of the L languages in the study, words in the K meaning categories are gathered. The Dyen et al. (1997) data uses the Swadesh (1952) "word list" (in fact a list of meanings). In this list, $K = 200$ core meaning categories ("All", "And", "Animal", ...) are identified. Words in the Swadesh meaning categories are relatively resistant to lateral trait transfer,

Table 1. A miniature lexical “data set” with $L = 3$ languages, $K = 3$ meanings and $N = 6$ distinct traits, $C = \{1, 2, \dots, 6\}$, from Dyen et al. (1997).

	“to give”	“big”	“we”			$k = 1$	$k = 2$	$k = 3$
Flemish	geven	groot	wy	\implies	$i = 1$	$c = 1$	$c = 3$	$c = 6$
Danish	give	stor	vi		$i = 2$	$c = 1$	$c = 4$	$c = 6$
Kashmiri	dyunu	bodü	asi		$i = 3$	$c = 2$	$c = 5$	$c = 6$

referred to here as *borrowing*. Embleton (1986) observes that words borrowed from French and Latin make up about 60% of the English lexicon, but less than 6% of the Swadesh 200-word list. The screened version of the Ringe et al. (2002) data we have uses a list of $K = 328$ meanings, (plus morphological traits, which we do not treat). There is a Swadesh list of $K = 100$ meaning categories thought to be particularly resistant to borrowing. The word lists are nested, so both data sets include the 200-word and 100-word lists.

In the following, trait data collected by Gray and Atkinson (2003) for Hittite, Tocharian A and Tocharian B are analysed with 84 languages (displayed in Fig. 4) from the Dyen et al. (1997) data. These merged data are referred to hereafter as the Dyen et al. (1997) data. Of the $L = 24$ languages in the Ringe et al. (2002) data, 20 are ancient. In contrast, of the $L = 87$ languages in the Dyen et al. (1997) data, just the three added by Gray and Atkinson (2003) are ancient. The two data sets are substantially independent. Both data sets are available on line at locations given in the bibliography.

The linguist identifies homology classes among the words in a given meaning category. In order to avoid false identification of homology, where there is merely a chance likeness of sound, linguists require close correspondence of meaning. Where words are judged to be descended from a common ancestor they are assigned the same trait label. This operation, which requires expert knowledge, is equivalent to replacing words with trait labels, $c \in C$, and thereby generating for each language $i = 1, 2, \dots, L$ and each meaning category $k = 1, 2, \dots, K$ a trait set $H_i^{(k)}$. In the context of this application, homology classes of traits are called *cognate classes*. The Gray and Atkinson (2003) registration of the Dyen et al. (1997) data which we treat here has $N = 2665$ cognate classes, and is presented to us as a 87×2665 binary matrix. The corresponding registration of the screened Ringe et al. (2002) data has $N = 3174$ cognate classes. In the example in Table 1, the data is coded $H_1^1 = \{1\}$, $H_2^1 = \{1\}$, $H_3^1 = \{2\}, \dots, H_3^3 = \{6\}$. Comparing with Fig. 1, we have here an extra superscript (k) on trait-sets H_i marking the meaning class subset. In Section 8 we take one independent copy of the trait birth-death process $H(\tau, i)$ for each meaning category.

Both data sets mark some homology classes as equivocal, and offer “splitting” and “lumping” versions of the data. We present results for the “splitting” data which assigns separate labels to homology classes which may in fact display a single homologous trait. Results for the lumping data are very similar. We comment on this systematic error in Section 9.5.

The vocabularies of some ancient languages are only partially reconstructed, creating gaps in the binary sequence data. The Ringe et al. (2002) data marks these gaps. We are unable to treat missing data at this stage. We are obliged to drop from the analysis of the Ringe et al. (2002) data the languages Gothic, Lycian, Luvian, Oscan, Umbrian, Old Prussian, Old Persian, Avestan and Tocharian A, leaving the languages in Fig. 8. We retain some languages with small numbers of gaps, simply marking the gap as trait-absence. We

discuss the associated model mis-specification bias in Section 9.6. The number of gaps in our registration of the Dyen et al. (1997) data is negligible.

Historical sources provide rate calibration data for these Indo-European data sets. Gray and Atkinson (2003) and Atkinson et al. (2005) compile calibration points. For example, the Celtic language-taxa *Irish_A*, *Irish_B*, *Welsh_N*, *Welsh_C*, *Breton_List*, *Breton_SE* and *Breton_ST* form a clade in the Dyen et al. (1997) data, with a common ancestor between 1700 and 2650 years before the present (BP, where the present is the year 2000 - only roughly the time the data was gathered, because the dating accuracy is in any case low). This is prior knowledge. Likewise, in the Ringe et al. (2002) data, *oldirish* and *welsh* form a clade with a common ancestor between 1700 and 2650 years BP. Analysis of the Dyen et al. (1997) data imposed 16 clades: Celtic, Brythonic, Italic, IberianFrench, Germanic, WestGermanic, NorthGermanic, BaltoSlav, Slav, Indic, IndoIranian, Iranian, Albanian, Greek, Armenian and Tocharic (see for example Figure 4). Analysis of the Ringe et al. (2002) data imposed the same prior constraints (wherever relevant).

Calibration points marked by bars in Figures 5 and 8 (bottom) are those that put lower and upper bounds on clade root times. Each such calibration point gives an independent estimate for λ , μ and θ . Prior knowledge which provides only a lower bound on language branching (“languages A and B were distinct by year C”) is more common, but less valuable, as it does not break the scale invariance discussed in Section 5.

8. Inference

8.1. Models

When we fit the model of Section 2 to the Dyen et al. (1997) and Ringe et al. (2002) data, we identify a model mis-specification problem. For meaning classes $k = 1, 2, \dots, K$ denote by $H^{(k)}(\tau, i)$ a trait birth-death process modelling the evolution of words in meaning category k , so that for $i = 1, 2, \dots, L$, $H_i^{(k)} = H^{(k)}(t_i, i)$ is the data at the leaves under NOABSENT. Let $\lambda^{(k)}$ and $\mu^{(k)}$ be the birth and death rates for traits in meaning class k . It is reasonable to expect any real language to have at least one word in each of the semantic fields in the Swadesh 200-word list at all times. It follows that the birth-death process must satisfy a *no-empty-field* condition, $H^{(k)}(\tau, i) \neq \emptyset$ or $N(\tau, i) > 0$, for each $(\tau, i) \in [g]$.

We ignore this no-empty-field condition in our analysis. This point is discussed further at the start of Section 9 and in Section 9.4. We lump together the K copies of the birth-death process of traits corresponding to the different meaning classes. Under the empty-field approximation, and assuming the death rates $\mu = \mu^{(k)}$, $k = 1, 2, \dots, K$ are all equal (see Section 9.3.1), the superposition

$$H(\tau, i) = \bigcup_{k=1}^K H^{(k)}(\tau, i) \quad (12)$$

of birth-death processes generates another instance of the same process, with birth rate $\lambda = \sum_k \lambda^{(k)}$ and death rate μ .

We carry out MCMC from posterior distributions $p(g, \mu | D)$ determined by Equation (9) and the Dyen et al. (1997) and Ringe et al. (2002) data, under the NOUNIQUE observation model. We repeated the analysis with NOABSENT for the Ringe et al. (2002) data, obtaining similar results. We apply the branching process prior f_G with hyperprior $1/\theta$, and uniform root prior f_R with $T = 16000$ (an uncontroversial upper limit on t_R). The

data overwhelm these two priors, differences between posterior estimates obtained under the two priors are slight, and we therefore present results for prior f_R exclusively. Results are completely insensitive to the choice of T , for all T sufficiently large.

In our search for conflicting signals in the data, we analyzed (in addition) subsets of the data. As discussed in Section 9.1, analyses of subsets of languages may be less exposed to error due to certain forms of borrowing. On the other hand we may uncover rate heterogeneity between word lists or between groups of languages. We reduce the Dyen et al. (1997) data to the Swadesh 100-word list, and the Ringe et al. (2002) data to the Swadesh 200- and 100-word lists. We thin the Dyen et al. (1997) data from $L = 87$ languages down to two sets containing $L = 31$ languages, and $L = 30$ languages, chosen in such a way that the pivotal calibrating dates remain applicable. These two data subsets overlap at 8 languages, but just one of the five calibration points has any common data (Tocharian, where there is no choice). We label analyses “Data/Prior/Word List/Leaves”.

8.2. Results

Figures 2 and 3 give a compact quantitative summary of the Dyen/ f_R /200/87, Dyen/ f_R /100/87, Dyen/ f_R /200/31, Dyen/ f_R /100/31, Dyen/ f_R /200/30, Dyen/ f_R /100/30 Ringe/ f_R /328/15, Ringe/ f_R /200/15 and Ringe/ f_R /100/17 posterior distributions. The posterior probabilities for a selection of clades are displayed in Fig. 2, and clade labels **BGCI**, **GCI**, **CI**, **CG**, **GI**, **GrA**, **noht** and **noht** defined. The posterior mean age for the common ancestor of the languages defining each corresponding clade is displayed in Fig. 3. This format is useful for identifying conflict between data subsets, once clades of interest have been identified.

The MCMC output was further processed using the R package Ape 1.8-2 of Paradis et al. (2004). We display consensus tree reconstructions of topology and branch length for the Dyen/ f_R /100/87 posterior (Fig. 4), the Dyen/ f_R /200/31 and Dyen/ f_R /100/30 posteriors (Fig. 6), and the Ringe/ f_R /100/17 posterior (Fig. 8 centre). A tree edge, or “split”, determines a partition of the leaves into two sets. Our majority rule consensus trees show all splits present in 50% or more posterior samples, and are treelike but multifurcating, and indicate the percent support for splits with support less than 95%. Edge length depicted is the posterior mean elapsed time in years, conditioned on the existence of the corresponding split. The consensus tree is a popular central point estimate in phylogenetic analyses. However, splits with posterior support close to 100% in a Bayesian phylogenetic analysis often have far lower frequency in bootstrap analyses of the same data. Here, the Dyen/ f_R /200/31 and Dyen/ f_R /100/30 consensus trees group Balto-Slav with Germanic, Celtic and Italic at above 95%, whilst Ringe/ f_R /100/17 puts Balto-Slav with Indo-Iranian at 90%.

The consensus tree is not a state in the sample space of trees. Prior constraints, and in particular leaf ages, are best represented on sampled states so we give, in Figure 5 samples drawn from the Dyen/ f_R /200/31, Dyen/ f_R /100/30 posterior distributions, and in Figures 8 and 9 samples from the Ringe/ f_R /100/17 and Ringe/ f_R /200/15 distributions.

Following Holland et al. (2005), we represent uncertainty in topology *via* a consensus network. The networks in Figures 7 and 8 were computed using the network analysis tool SplitsTree V4.4 of Huson and Bryant (2006). They display all splits with posterior support above 5%. The consensus network contains the consensus tree as a subtree. Branch lengths measure posterior support. We know of no analysis or software treating rooted consensus networks.

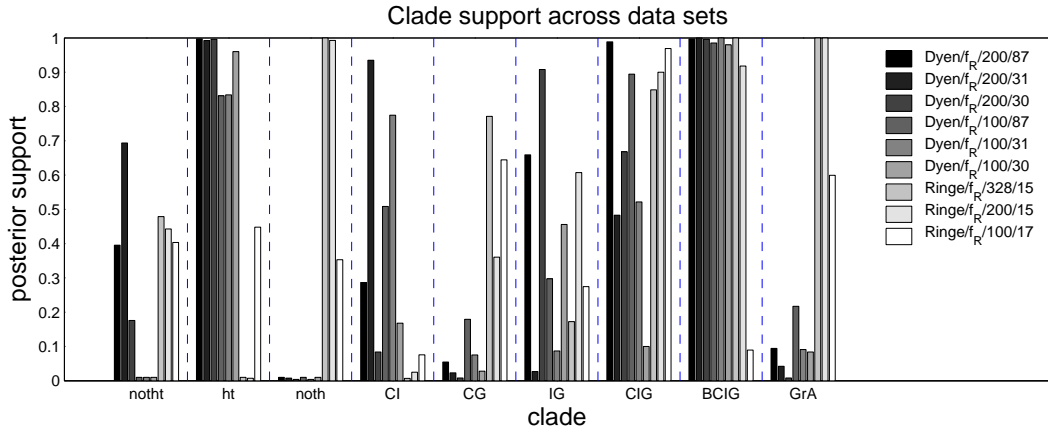


Fig. 2. Posterior probabilities for selected clades, across data sets. *x*-axis labels: BGCI, Balto-Slav-Germanic-Celtic-Italic; GCI, Germanic-Celtic-Italic; CI, Celtic-Italic; CG, Celtic-Germanic; GI, Germanic-Italic; GrA, Greek-Armenian; noht, complement of Hittite-Tocharian; noth, complement of Hittite.

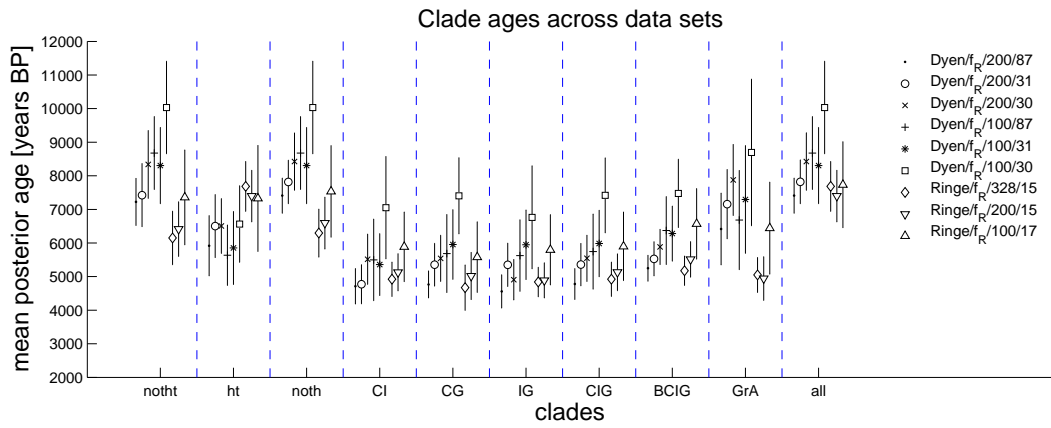


Fig. 3. The mean posterior ages, in years BP, for the most recent common ancestor (MRCA, super-clade root) of languages in selected combinations of clades (*ie* super-clades). *x*-axis labels as for Fig. 2.

8.3. Discussion

There is some conflict in the support for clades across analyses. For example, while the **BCIG** group is strongly supported in all analyses (except Ringe/ f_R /100/17, which allows it, as we see from the network in Fig. 8, and the bar in Fig. 2) and although the **CIG** group is supported in all analyses (except Dyen/ f_R /100/30, which allows it, as we see from the bar in Fig. 2 - the split is present in a tangle in the network in Fig. 6), all the sub-clades **CI**, **CG** and **IG** conflict at least one analysis. Age estimates for the common ancestor of the languages in the **CIG** clade are, in all analyses, close to the age estimates for the common ancestors of the subclades, suggesting the breakup occurred in a relatively small interval of time, so the split structure is poorly resolved. Conflicting support for topologies may be due to borrowing in the breakup or subsequent evolution. Age estimates are robust to this form of model misspecification. See Section 9.1.

Referring to the consensus trees, Figures 6 and 8, and the clade probabilities, Fig. 2, notice that **HT** (Hittite and Tocharian) is an outgroup in the three Dyen/ f_R /200/ Y analyses, grouped with Greek and Armenian in the Dyen/ f_R /100/ Y analyses, and split in the three Ringe/ f_R / Y/X analyses. There are many model misspecification issues for Hittite and Tocharian. Comparing **noth** and **all** in Fig. 3, Hittite adds a thousand years to the posterior mean root age of the Ringe/ f_R /328/15 and Ringe/ f_R /200/15 analyses. The contrast between the Dyen et al. (1997) and Ringe et al. (2002) analyses is most clearly visible in the **noth** and **notht** columns of Figures 2 and 3.

In contrast, conflicts between analyses of the Swadesh 100-word list (Dyen/ f_R /100/87, Dyen/ f_R /100/31, Dyen/ f_R /100/30 and Ringe/ f_R /100/17) are almost absent. Both **CI** (visible in Fig. 8) and **IG** (see Fig. 2) are allowed by these analyses. Ringe/ f_R /100/17 allows the clade **HT** and does not impose **HT** as an outgroup (which would be a conflict, as **notht** is not a clade of Dyen/ f_R /100/ Y). In other areas of conflict, the Dyen/ f_R /100/ Y analyses allow a **GrA** clade. This lack of conflict comes at the price of greater random error (compared to analyses on longer word-lists). One striking conflict remains: the position of Indo-Iranian relative to the root is quite different in the Ringe/ f_R /100/17, and Dyen/ f_R /100/ Y analyses.

The posterior mean ages for the **noth**, **notht**, **all** and **GrA**, which show particular conflict in Figures 2 and 3, are in agreement for analyses based on the Swadesh 100-word list. This reduced set of traits is chosen to be resistant to borrowing. Posterior predictive replicates computed in Section 9.3.1 show little evidence of rate heterogeneity within this class of traits. The corresponding words are relatively well attested in otherwise incompletely reconstructed ancient languages, so there is little missing data. In Section 9.2 we compute posterior predictive distributions for unique traits in the Ringe/ f_R /100/17 analysis; these agree well with external data.

In summary, the systematic errors displayed in our four age estimates from the Swadesh 100-word list are representative. On the other hand, most features of tree topology which were doubtful, remain in doubt.

9. Model mis-specification

We head this section with a summary of its results. These results coincide with the conclusions we draw from the between-data analyses in Section 8.3: our age estimates are robust; tree topology less so. In Fig. 10 and Fig. 11 we present results from synthetic data, simulated on the tree in Fig. 9, under a range of observation models intended to mimic likely model mis-specification. Details of these models, which simulate the empty-field-condition,

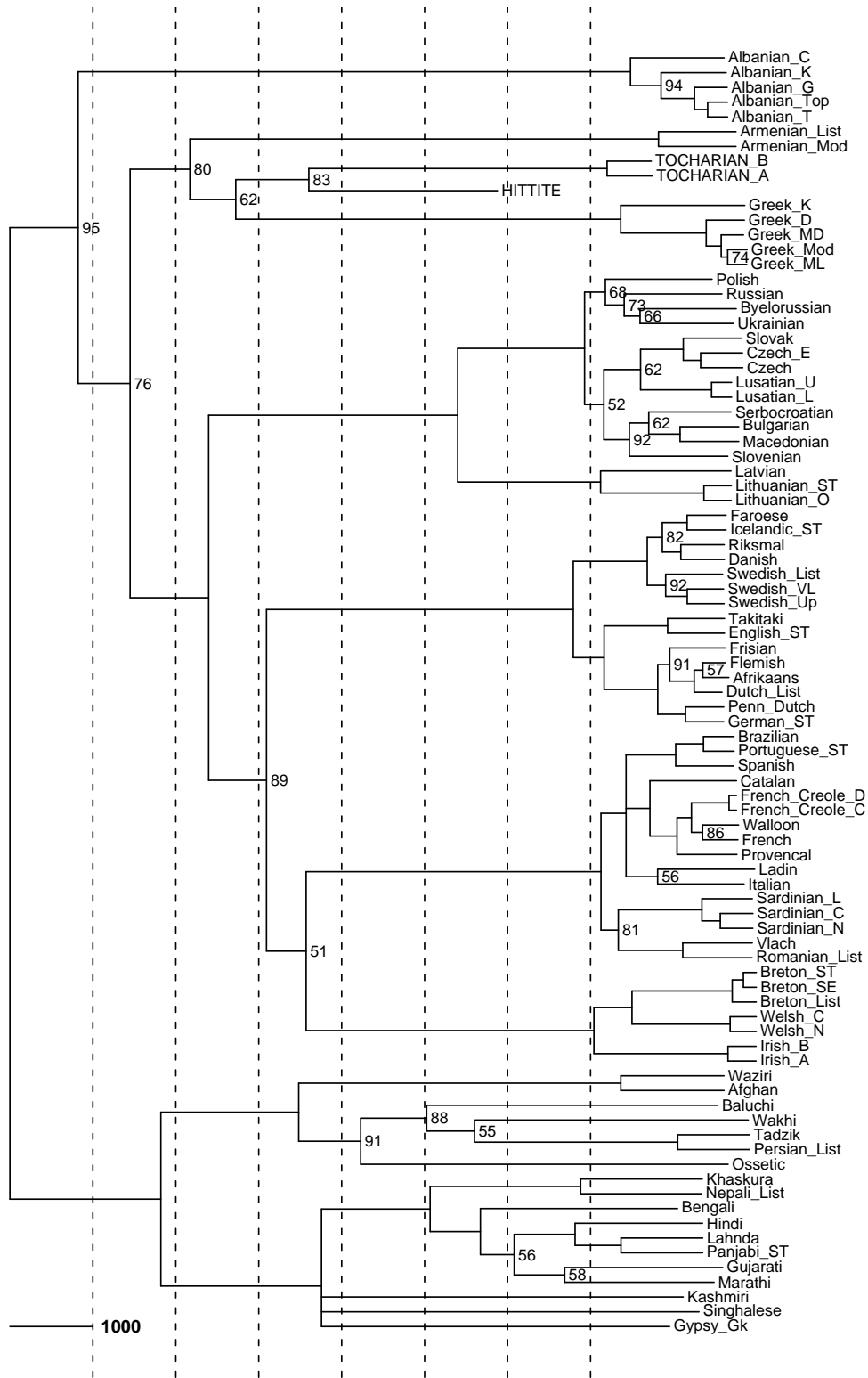


Fig. 4. Consensus tree for the $D_{yen}/f_R/100/87$ posterior distribution. Labelled nodes are supported at less than 95% posterior probability.

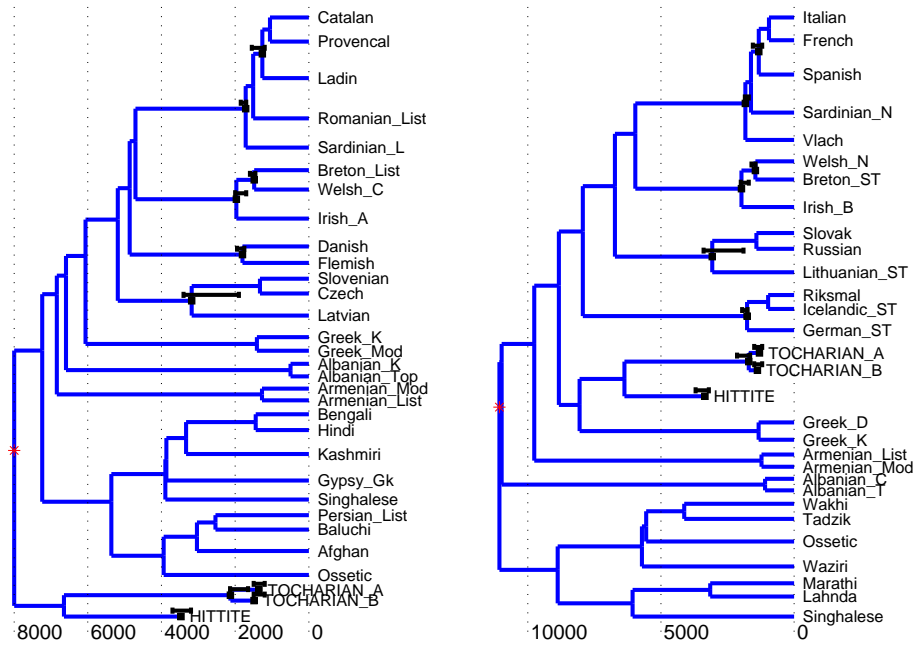


Fig. 5. Trees sampled from the $D_{\text{yen}}/f_R/200/31$ (left) and $D_{\text{yen}}/f_R/100/30$ (right) posterior distributions. x -axis gives age in years. Prior constraints on seven clade root and three leaf ages are indicated by vertical error bars.

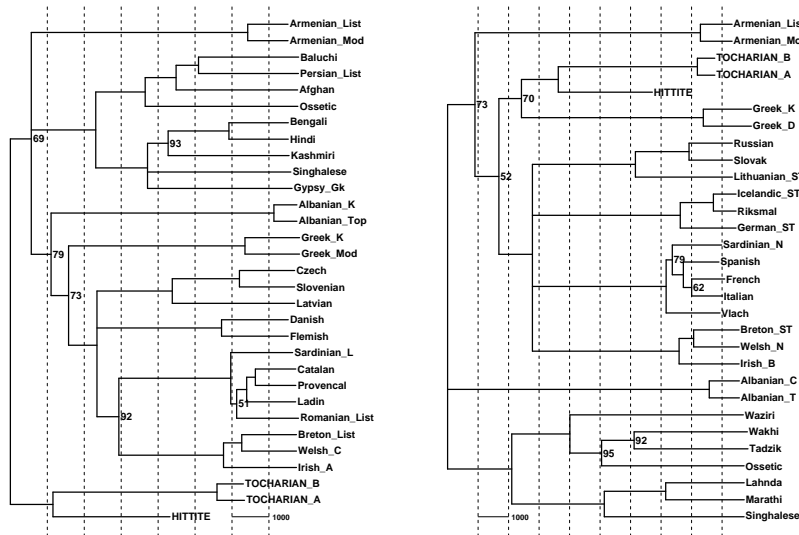


Fig. 6. Consensus trees computed for the $D_{\text{yen}}/f_R/200/31$ (left) and $D_{\text{yen}}/f_R/100/30$ (right) posterior distributions.

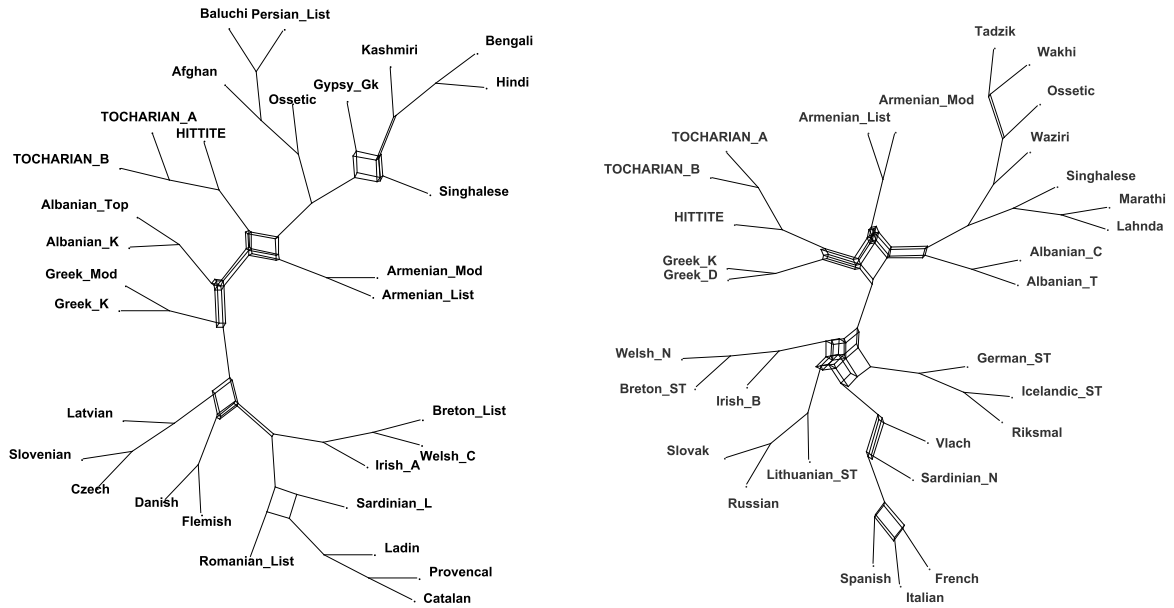


Fig. 7. Consensus networks summarising the Dyen/ $f_R/200/31$ (left) and Dyen/ $f_R/100/30$ (right) posterior distributions. Branch length is proportional to posterior support for the corresponding split.

plausible levels of borrowing and branch-wise and trait-wise rate heterogeneity, are given in Sections 9.1 through 9.6.

Clades imposed in reconstructions from synthetic data are indexed **S-** and are the same as the clades defined below Fig. 2. Analyses of synthetic data are indexed “S/X/Y”. Values of X indicate borrowing and rate heterogeneity: X=“T” is no borrowing; X=“Gb” is global borrowing at rate $b\mu$; X=“Lz - b” is local borrowing, between languages with a common ancestor not more than z years in the past, at rate $b\mu$; and X=“BH ρ ” and X=“MH ρ ” have rates drawn independently, for each branch (BH) and meaning category (MH) respectively, from a Gamma distribution with mean μ and standard deviation $(\rho/100)\mu$. Values of Y show the constraint applied: Y=“Un” is the unconstrained birth death process of set elements, with $n = \lambda/\mu$ the expected number of distinct traits at each leaf, under the NOABSENT observation model; Y=“Cn” simulates cognate classes under the no-empty-field constraint, using n meaning categories.

Systematic errors in tree *node ages* inferred from synthetic data generated under these models are in general small. With the exception of $S/BH50/U200$ (see Section 9.3.2), the systematic error we generated in Fig. 11 is of the same order of magnitude as, or smaller than, random error. Systematic error in estimated *rates* μ (not shown) is highly significant. Calibration data fixes dates and topology for that part of the tree adjacent to the leaves, forcing the inference to accommodate model mis-specification by adjusting rates. The modified rates fit the imposed trait evolution adjacent to the leaves. If model mis-specification is homogeneous over the tree, as is the case for the empty-field-approximation, trait evolution deep in the tree may be well represented by these biased rates, and date

estimates are accordingly robust.

Tree topology is not robust to the model mis-specification we explored. The “true” s-clades **s-GrA**, **s-BCIG**, and the rooting clades **s-notht**, and **s-noth** have robust support at levels which do not allow rejection of the truth. The “true” s-clades **s-CIG** and **s-IG** are not well reconstructed when borrowing is substantial. The branching at the top of the superclade **s-BCIG** in Fig. 9, centre, is poorly resolved as the **s-BCIG**, **s-CIG** and **s-IG** branches are separated by just 1000 years, which is small compared to $\mu^{-1} \simeq 5000$. Nevertheless, the truth is rejected only at very high levels of borrowing (S/Gb/Y where $b = 0.2, 0.5$). Clade age estimates, shown in Fig. 3 and Fig. 11, can be stable across analyses when topology is uncertain. This is because the super-clade ages are determined largely by total tree length; total tree length is tightly coupled to the number of transitions on the tree, which is rather well determined by data.

9.1. Global and local borrowing

Word borrowing is akin to lateral gene transfer on the tree of life. Word borrowing from languages outside the study is straightforward trait birth (unless the same word is borrowed into several languages). Moreover, if we delete a source language from our study, we remove the model error associated with borrowing from that language. The consistency we see in Fig. 3 between clade ages reconstructed for near-disjoint subsets of languages, and the full set, in the Dyen/ $f_R/200/87$, Dyen/ $f_R/200/31$ and Dyen/ $f_R/200/30$ posteriors, suggests that borrowing is not distorting the Dyen/ $f_R/200/87$ estimates themselves.

Our models of borrowing are as follows. We associate with each time slice $\tau \in [0, \infty)$ across the tree a linkage graph $(E(\tau), V(\tau))$ with nodes, $V(\tau) = \{(\tau, i) \in [g]\}$, corresponding to points in $[g]$ intersected by the slice. The linkage graph models traffic between languages; its edges $\langle y, z \rangle \in E(\tau)$ connect sets $H(y)$ and $H(z)$ between which trait instances can pass. Let $\tilde{V}(z) = \{z \in V(\tau) : \exists y \in V(\tau), \langle y, z \rangle \in E(\tau)\}$ be the set of nodes adjacent to $z \in V(\tau)$. Let b denote the relative rate of word-borrowing to word-death. At per capita rate $b\mu$ each instance of each trait in the time slice τ generates a borrowing event. Suppose the selected trait-instance is labelled $c \in C$. A language $y \in \tilde{V}(z)$ is chosen, uniformly at random from nodes adjacent to y on $(E(\tau), V(\tau))$, and we set $H(y) \leftarrow H(y) \cup \{c\}$, *ie* the word is copied into the target language.

We model local borrowing as follows. Words transfer between languages which have a sufficiently recent common ancestor. The linkage graph at time t includes an edge from (t, i) to (t, j) if points (t, i) and (t, j) in $[g]$ have a common ancestor less than z years in the past. In this model linked groups of languages break up into linked subgroups. In our model of widespread borrowing (the “global” borrowing model), all languages communicate equally with all other languages, and the linkage graph is the complete graph.

Our exploration of these models is summarized in Figures 10 and 11 by the three S/Gb/U200 data sets and the S/L500 – 1/U200 data set. We display global borrowing at relative rates of 10%, 20% and 50% the death rate. Higher global rates are probably irrelevant: our own studies of synthetic data show that very high borrowing rates lead to inconsistency between analyses of the full dataset and subsets of languages. No such inconsistency is seen. Local borrowing has time depth $z = 500$ and a borrowing rate equal to the death rate. We see from Fig. 11 that age estimates are robust to this form of model mis-specification.

9.2. Predictive distributions and external data

Where the observation model is NOABSENT, unique traits *are* present, and we can use them to test the model. We drop them from the data, carry out the inference under NOUNIQUE, and then see if we can predict the number of unique traits for each taxon. This check was available for the Ringe et al. (2002) data. We expect rate heterogeneity and borrowing to be visible (but probably not distinguishable) in these tests. Results are shown in Fig. 12.

Denote by \tilde{D} synthetic trait data generated under the NOABSENT observation model, displaying \tilde{N} distinct traits. For trait $a = 1, 2, \dots, \tilde{N}$ let \tilde{M}_a give the indices of leaves displaying an instance of trait a for predicted data \tilde{D} , and \tilde{X}_i be the number of unique traits in \tilde{D} at taxon i ,

$$\tilde{X}_i = \text{card} \{ \tilde{M}_a : \tilde{M}_a = \{i\}, a = 1, 2, \dots, \tilde{N} \} \quad i \in V_L.$$

The posterior predictive distribution $\Pr\{\tilde{D}|D\}$ is

$$\Pr\{\tilde{D}|D\} = \int \Pr\{\tilde{D}|g, \mu, \lambda\} p(g, \mu, \lambda|D) dg d\mu d\lambda$$

and this determines a predictive distribution for \tilde{X}_i . We sample μ, λ and g from the posterior $p(g, \mu, \lambda|D)$ (g and μ are available from MCMC output; we restore λ by sampling its posterior conditional density), simulate synthetic data \tilde{D} at the leaves of g , and compute \tilde{X}_i from \tilde{D} . Let $X_i(D) = \text{card} \{ M_a : M_a = \{i\}, a = 1, 2, \dots, N \}$, denote the number of unique traits at taxon i in the data itself.

Predictive distributions for X_i are shown in Fig. 12. The predictive distributions for \tilde{X}_i over-estimate the X_i in the full Ringe et al. (2002) data. Since borrowing depletes unique traits, this is consistent with model mis-specification due to borrowing. Also, we expect borrowing to be weaker on the shorter word-lists, since the shorter lists are by design more resistant to borrowing. We see in Fig. 12 that unique traits are indeed more reliably predicted on shorter lists. Corresponding studies are given for synthetic data in Fig. 13, simulated with decreasing levels of borrowing, and show a pattern of decreasing over-estimation with decreasing borrowing in a synthetic vocabulary of fixed size. The posterior predictive distributions presented here are not specifically a test for borrowing. They are, rather, a health check on the data and model. Predictive distributions from the shorter word-lists are in good agreement with the data.

9.3. Rate heterogeneity

Our estimates of systematic error in reconstructed dates due to rate heterogeneity are based on the rate variation shown in Fig. 14. We measure the per-trait-instance death rate μ for each clade calibration constraint independently. We sampled the posterior distributions $p(g, \mu|D_{\text{clade}})$ determined by the data for each clade in turn.

9.3.1. Rate heterogeneity across traits

Pagel and Meade (2006) show that the evolution rates of words are, for a given meaning category, fairly consistent across data sets, whilst varying more substantially between meaning categories. Fig. 14 displays a tendency for the shorter word-lists to evolve at relatively slower rates. This is expected. However, comparing Figures 3 and 14, we see that rate variation between data sets does not lead directly to variation in estimated root times. For

example, Dyen/ $f_R/200/87$ and Dyen/ $f_R/100/87$ differ by a factor 1.5 in posterior mean rate, but by just 1.2 in root age. Time depth measurements do not depend on an assumption of constant rates *between* analyses, since rates are estimated from calibration points in the recent history of the same data used to predict branching times.

In order to generate synthetic data with rate heterogeneity across meaning classes (the $S/MH\rho/U200$ simulations), we draw rates

$$\mu^{(k)} \sim \text{Gamma}(\alpha, \beta)$$

independently for each meaning category $k = 1, 2, \dots, K$, with mean $\alpha\beta = \mu$ and variance $\alpha\beta^2 = (r\mu)^2$, where $r = 0.25, 0.5$ and $r = 1$, simulate a trait process $H^{(k)}(\tau, i)$ at rate $\lambda, \mu^{(k)}$, merge meaning categories, as Equation (12), and then read off data, as Equation (2). Pagel and Meade (2006) estimate rate variance over meaning classes $(r\mu)^2 \simeq \mu^2/9$.

Rate heterogeneity across traits distorts the distribution of the number, $\text{card } M_a$, of languages in which trait $a \in C$ appears. Denote by $Y^{(n)} = Y^{(n)}(D)$ the number of traits displayed at n leaves,

$$Y^{(n)}(D) = \text{card} \{M_a : \text{card } M_a = n, a = 1, 2, \dots, N\},$$

and let $\tilde{Y}^{(n)} = Y^{(n)}(\tilde{D})$ be the corresponding random variable computed from posterior predictive data $\tilde{D} \sim \text{Pr}\{\tilde{D}|D\}$. We plot $E(\tilde{Y}^{(n)}|D) - Y^{(n)}$ and the envelope $\pm 2\text{std}(\tilde{Y}^{(n)}|D)$. In Fig. 15 the model is unable to reproduce the trait frequency distribution in synthetic data with high levels of rate heterogeneity across meaning classes (standard deviation 50% of the mean) but lower levels (25%) are invisible. In the terminology of O'Hagen and Forster (2004) chapter 8, this posterior predictive replicate diagnostic is simply a consistency check, as opposed to Section 9.2, where we have external data.

Returning to the real data, in Fig. 16 we see some inconsistency attributable to rate heterogeneity between traits in the Ringe/ $f_R/328/15$ analysis. Among other problems, the data contains an excess of traits appearing in 10 or more leaves. This is caused by a small cohort of traits evolving at death rate μ small compared to the rest. The effect is very greatly reduced in the Ringe/ $f_R/100/17$ analysis at right. Analyses of the Dyen et al. (1997) data show a similar pattern.

9.3.2. Rate heterogeneity in space and time

Time depth measurements do depend on some assumption about the way rates have changed over time *within* each data set we analyze. The variations in rates between clades within each of the D100, D200, R328 and R100 groups in Fig. 14 give us an indication of the un-modelled rate variation we can expect in the deeper branches of the tree.

Synthetic data with spatio-temporal rate variation ($S/BH\rho/U200$ analyses), have rates

$$\mu_{\langle i,j \rangle} \sim \text{Gamma}(\alpha, \beta)$$

drawn independently on each edge $\langle i, j \rangle \in E$, with mean $\alpha\beta = \mu$ and variance $\alpha\beta^2 = (r\mu)^2$, where $r = \rho/100$, so that the standard deviation of the rates is 10%, 33% and 50% the mean rate. Lees (1953) sees 20% variation in rate estimates from pairs of languages. Results are robust to moderate levels of unstructured, random rate variation of this kind.

Results are of course not robust to structured rate variation, in which, for example, rates on edges at ages greater than any calibration point are all larger than any rates in

the calibration zone, or a single-taxon outgroup has an extreme rate. In *S/BH50/U200*, *s – hittite* happens to have a high rate. Its root is pushed to great age, and with it goes the root of the tree. The analysis is exposed to rare “catastrophic” trait-evolution events outside the calibration zone. Examination of the *Ringe/f_R/100/17*, and *Dyen/f_R/100/Y* consensus trees and clade probabilities show that no single language or small outgroup is determining the root age. Agreement between reconstructions based on the predominantly ancient languages of Ringe et al. (2002) and modern languages of Dyen et al. (1992) shows that there is at least no such structured rate variation in the recent past. It is feasible also to fit a model with explicit rate heterogeneity, as is common practice in phylogenetics.

9.4. The empty-field approximation

Our empty-field approximation will be good if there is significant “polymorphism”, that is, if the mean number $\lambda^{(k)}/\mu^{(k)}$ of traits (*ie*, words per meaning category) in the $H^{(k)}(\tau, i)$ -process is large. We estimate λ/μ at 273(9) for the Dyen et al. (1997) Swadesh-200 data and 280(25) for the Ringe et al. (2002) Swadesh-200 data (posterior standard deviation in parenthesis) and hence $\lambda^{(k)}/\mu^{(k)} \simeq 1.4$. The probability, $\exp(-\lambda^{(k)}/\mu^{(k)}) \simeq 1/4$, to find the unconstrained trait-set process $H^{(k)}(\tau, i)$ in the empty set at any single fixed point $(\tau, i) \in [g]$ is high enough to cause concern.

We simulate synthetic data from the trait birth-death process constrained to respect the no-empty-field condition. For each of the $k = 1, 2, 3, \dots, K$ meaning classes, we simulate $N^{(k)}(t_R, R)$ from a Poisson distribution constrained to be greater than zero, then simulate $H^{(k)}(\tau, i) | N^{(k)}(\tau, i) > 0$ in $[g]$. The total rate for the exponential waiting time to the next event does not include μ if $N^{(k)}(\tau, i) = 1$. We then merge the meaning classes as in Equation (12). Our studies are represented here by two simulations, *S/T/C200* and *S/L500-1/C200*, the latter including local borrowing. The per-capita death rate μ was set to a large value, so that polymorphism was low. We find, when we fit data of this kind, that the tree and its dates are robust to this form of model mis-specification.

9.5. Incorrect splitting deep rooted homology classes

When the scientist groups instances of traits into homology classes, instances of traits born deep in the tree may be highly evolved, and correspondingly difficult to identify as in fact homologous. This error can populate the deeper branches of the tree with spurious birth events. This is a case where model mis-specification is not homogeneous over the tree, and will lead to over-estimation of the tree depth. When we replace the Ringe et al. (2002) “splitting” data with the Ringe et al. (2002) “lumping” data we do see a 3% downward shift in the estimated root time.

9.6. Unknown vocabulary as absent traits

In our analysis of the Ringe et al. (2002) data, we retain some languages with gaps, corresponding to missing data. We replace these gaps with zeros, marking trait absence. Gappy languages (*Hittite*, *Tocharian*) do stand out in predictive tests counting singleton states on external data for the 328 and 200 word-lists. However, the effect is removed when we reduce the data to the Swadesh 100 word-list, where traits are better attested. The effect is to bias reconstructed branching times for gappy taxa to larger age values on the Ringe et al. (2002) 328 and 200 word-list data (see for example *ht* in Fig. 3).

Acknowledgements

The authors acknowledge advice and assistance from Quentin Atkinson and David Welch of the University of Auckland, and financial support from the Royal Society of New Zealand.

References

- Atkinson, Q., G. Nicholls, D. Welch, and R. Gray (2005). From words to dates: water into wine, magic or phylogenetic inference. *Transactions of the Philological Society* 103, 193–219.
- Bergsland, K. and H. Vogt (1962). On the validity of glottochronology. *Current Anthropology* 3, 115–153.
- Blust, R. (2000). Why lexicostatistics doesn't work: the "universal constant" hypothesis and the austronesian languages. In C. Renfrew, A. McMahon, and L. Trask (Eds.), *Time depth in historical linguistics*, pp. 311–332. Cambridge, UK: The McDonald Institute for Archaeological Research.
- Drummond, A. J., G. K. Nicholls, A. G. Rodrigo, and W. Solomon (2002). Estimating mutation parameters, population history and genealogy simultaneously from temporally spaced sequence data. *Genetics* 161, 1307–1320.
- Dyen, I., J. Kruskal, and B. Black (1992). An Indo-European classification: a lexicostatistical experiment. *Transactions of the American Philosophical Society* 82(1-132).
- Dyen, I., J. Kruskal, and P. Black (1997). FILE IE-DATA1, raw data available from <http://www.ntu.edu.au/education/langs/ielex/IE-DATA1>. Binary data available from <http://www.psych.auckland.ac.nz/psych/research/RusselsData.htm>.
- Embleton, S. (1986). *Statistics in Historical Linguistics*. Bochum: Brockmeyer.
- Erdem, E., V. Lifschitz, and D. Ringe (2005). Temporal phylogenetic networks and logic programming. *Theory and Practice of Logic Programming*. To appear.
- Felsenstein, J. (1981). Evolutionary trees from DNA sequences: a maximum likelihood approach. *J. Mol. Evol.* 17, 368–376.
- Felsenstein, J. (1992). Phylogenies from restriction sites: a maximum-likelihood approach. *Evolution* 46, 159–173.
- Geyer, C. J. (1992). Practical Markov chain Monte Carlo (with discussion). *Statist. Sci.* 7, 473–511.
- Gray, R. and Q. Atkinson (2003). Language-tree divergence times support the Anatolian theory of Indo-European origin. *Nature* 426, 435–439.
- Holland, B. R., F. Delsuc, and V. Moulton (2005). Visualizing conflicting evolutionary hypotheses in large collections of trees: Using consensus networks to study the origins of placentals and hexapods. *Systematic Biology* 54, 66–76.
- Huson, D. and D. Bryant (2006). Application of phylogenetic networks in evolutionary studies. *Molecular Biology and Evolution* 23, 254–267.

- Huson, D. and M. Steel (2004). Phylogenetic trees based on gene content. *Bioinformatics* 20(13), 2044–2049.
- Kimura, M. and J. Crow (1964). The number of alleles that can be maintained in a finite population. *Genetics* 49, 725–738.
- Lees, R. (1953). On the basis of glottochronology. *Language* 29, 113–127.
- Nakhleh, L., D. Ringe, and T. Warnow (2005). Perfect phylogenetic networks: A new methodology for reconstructing the evolutionary history of natural languages. *LANGUAGE, Journal of the Linguistic Society of America* 81, 382–420.
- O’Hagen, A. and J. Forster (2004). *Bayesian Inference* (2nd ed.), Volume 2B of *Kendall’s Advanced Theory of Statistics*. Arnold.
- Pagel, M. and A. Meade (2006). Estimating rates of lexical replacement on phylogenetic trees of languages. In P. Forster and C. Renfrew (Eds.), *Phylogenetic Methods and the Prehistory of Languages*, pp. 173–181. Cambridge, UK: The McDonald Institute for Archaeological Research.
- Paradis, E., J. Claude, and K. Strimmer (2004). Ape: Analysis of Phylogenetics and Evolution in R language. *BIOINFORMATICS* 20, 289–290.
- Ringe, D., T. Warnow, and A. Taylor (2002). Indo-European and Computational Cladistics. *Transactions of the Philological Society* 100, 59–129. Data available from <http://www.cs.rice.edu/~nakhleh/CPHL/#datasets>.
- Ronquist, F. and J. P. Huelsenbeck (2003). Mrbayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics* 19, 1572–1574.
- Saitou, N. and M. Nei (1987). The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution* 4, 406–425.
- Sankoff, D. (1973). Mathematical developments in lexicostatistic theory. *Current Trends in Linguistics* 11, 93–113.
- Swadesh, M. (1952). Lexico-statistic dating of prehistoric ethnic contacts. *Proc. Am. Phil. Soc.* 96, 453–463.
- Warnow, T., S. Evans, D. Ringe, and L. Nakhleh (2006). A stochastic model of language evolution that incorporates homoplasy and borrowing. In P. Forster and C. Renfrew (Eds.), *Phylogenetic Methods and the Prehistory of Languages*, pp. 75–87. Cambridge, UK: The McDonald Institute for Archaeological Research.
- Watterson, G. (1975). On the number of segregating sites in genetical models without recombination. *Theoretical Population Biology* 7, 256–276.
- Yang, Z. and B. Rannala (1997). Bayesian phylogenetic inference using dna sequences: a markov chain monte carlo method. *Molecular Biology and Evolution* 14, 717–724.

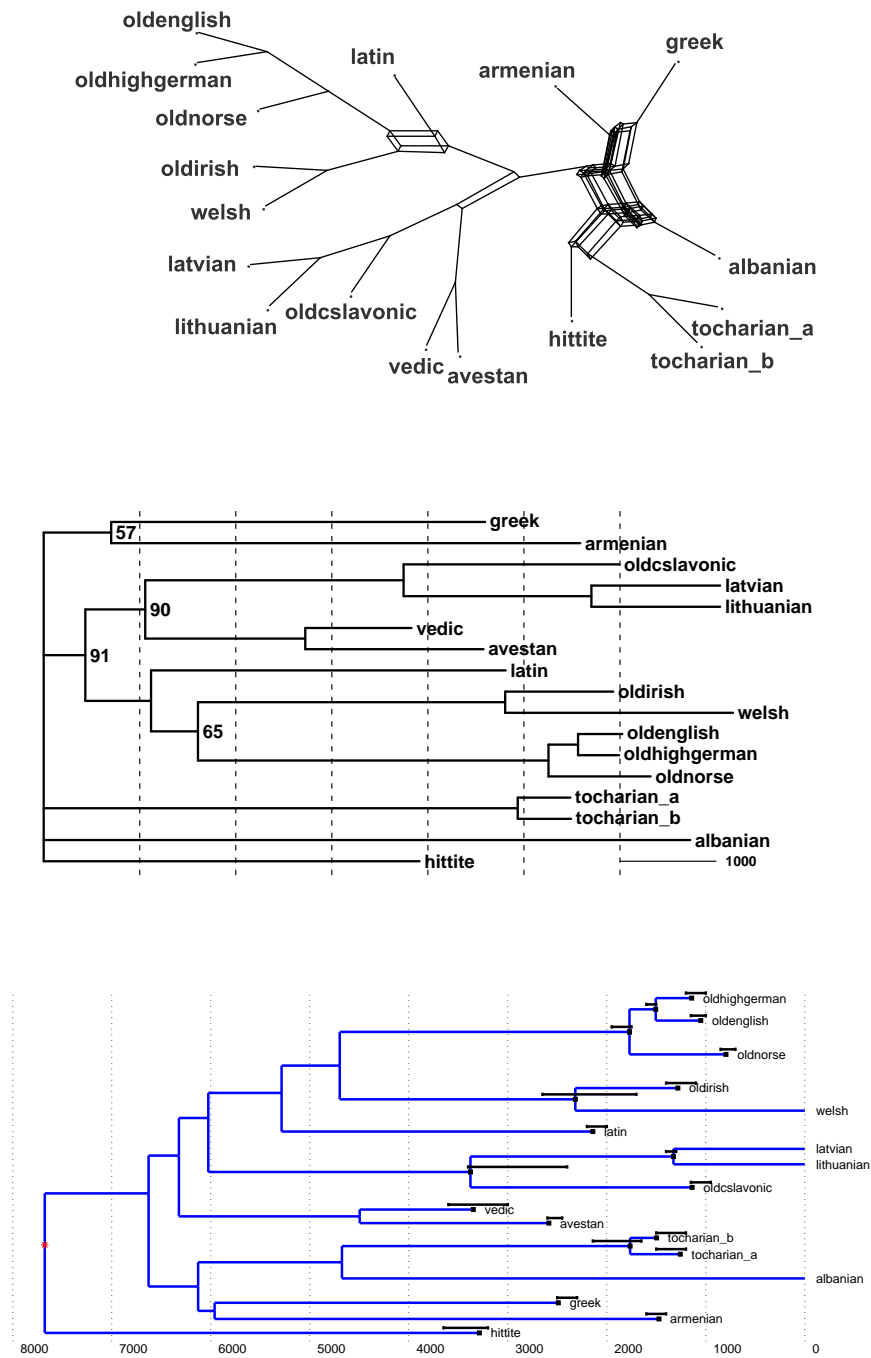


Fig. 8. Consensus network (top), consensus tree (centre) and a sampled state (bottom) illustrating the $\text{Ringe}/f_R/100/17$ posterior distribution. Consensus network and tree as Figures 7 and 6. Prior constraints on topology impose 7 clades. Prior uncertainties in clade root and leaf ages are indicated by vertical error bars (bottom).

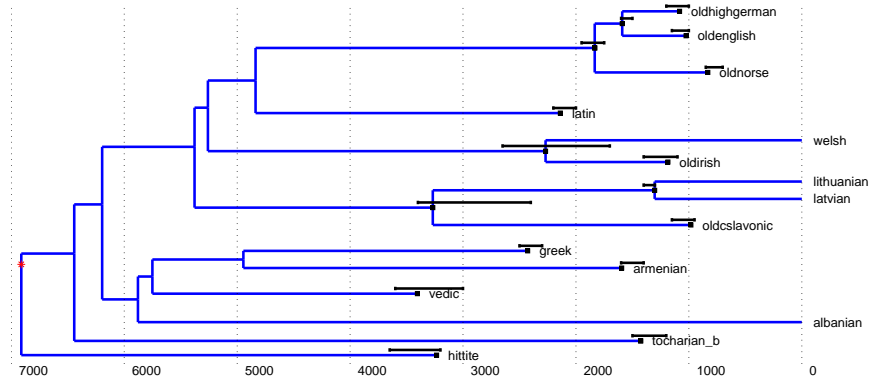


Fig. 9. A tree sampled from the $\text{Ringe}/f_R/200/15$ posterior distribution, the truth in the synthetic studies reported below.

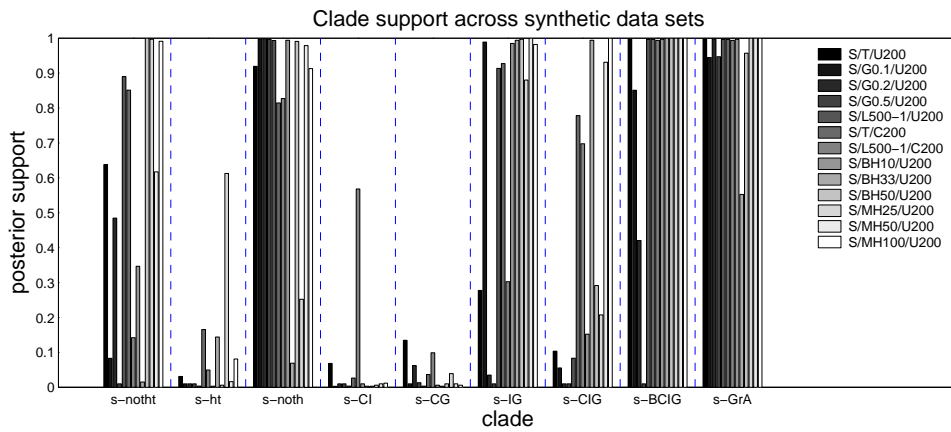


Fig. 10. Synthetic data yields estimates of posterior probabilities for selected clades, across synthetic data sets. x -axis labels as for Fig. 2 with s- prefix indicating synthetic.

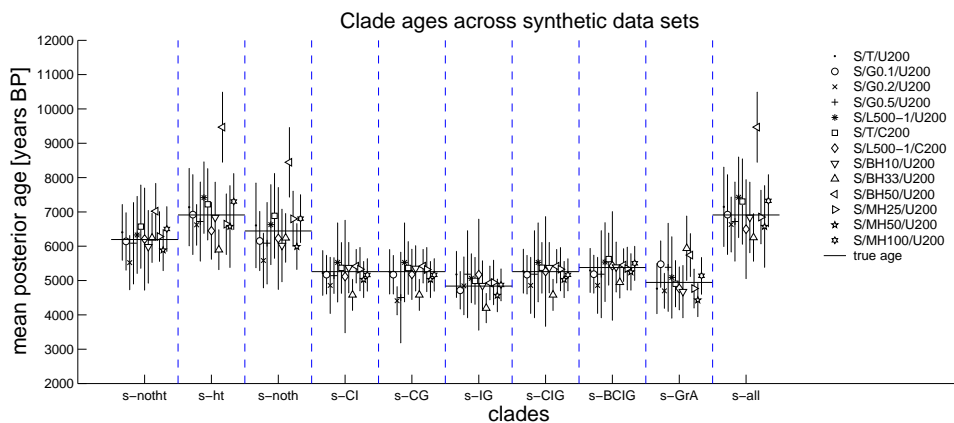


Fig. 11. Synthetic data yields estimates of mean posterior ages, in years BP, for the most recent common ancestor (MRCA, super-clade root) as in Fig. 2 with s- prefix indicating synthetic.

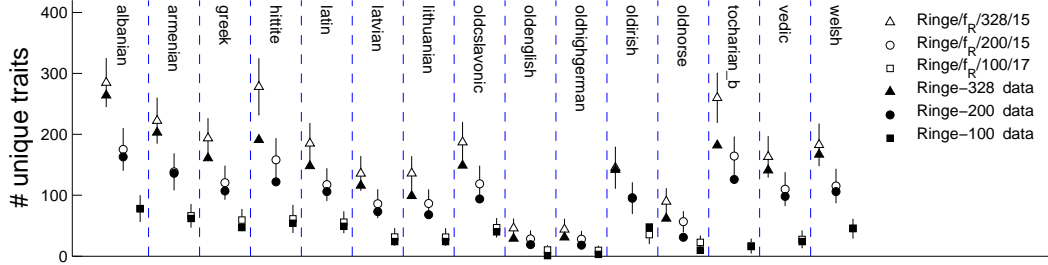


Fig. 12. Posterior predictive distributions for the number of unique states, external data from the Ringe et al. (2002): \triangle 328 meaning categories; \circ Swadesh-200 list; \square Swadesh 100 list.

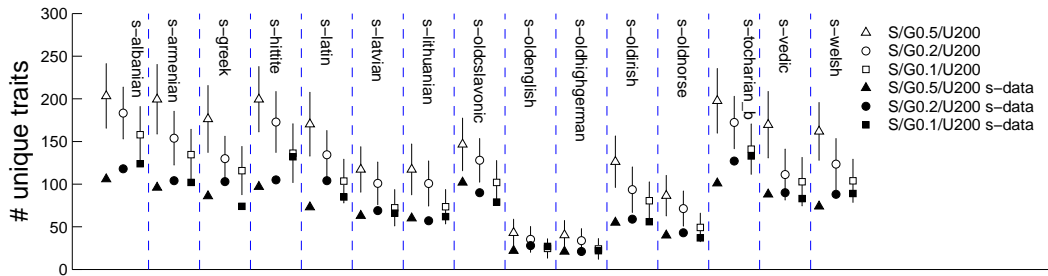


Fig. 13. Synthetic data incorporating borrowing yield posterior predictive distributions for the number of unique states, and show increasing unique-trait depletion as the level of synthesized borrowing increases: \triangle borrowing rate $\beta = 0.5\mu$; \circ borrowing rate $\beta = 0.2\mu$; \square borrowing rate $\beta = 0.1\mu$.

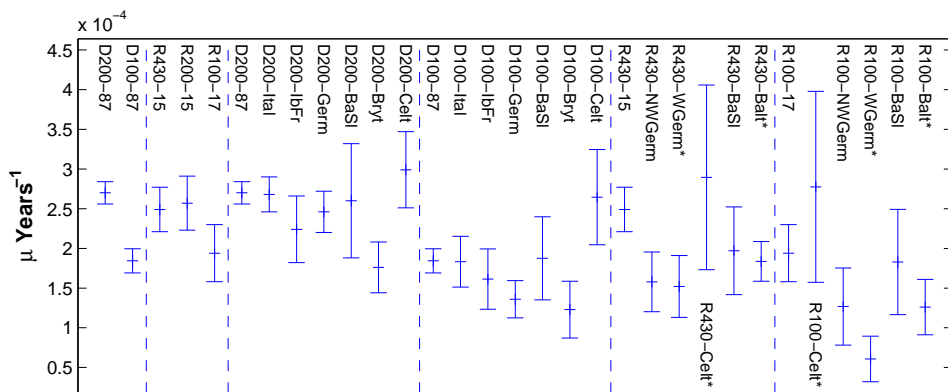


Fig. 14. Posterior mean values of μ , with standard errors at two standard deviations, measured in analyses Dyen/ f_R /200/87 (D200-87), Dyen/ f_R /100/87 (D100-87), Ringe/ f_R /328/15 (R328-15), Ringe/ f_R /200/15 (R200-15) and Dyen/ f_R /100/17 (R100-17), and independently using calibration points in distinct clades. The observation model is NOUNIQUE, except for two-leaf clades marked *, where NOABSENT must be used.

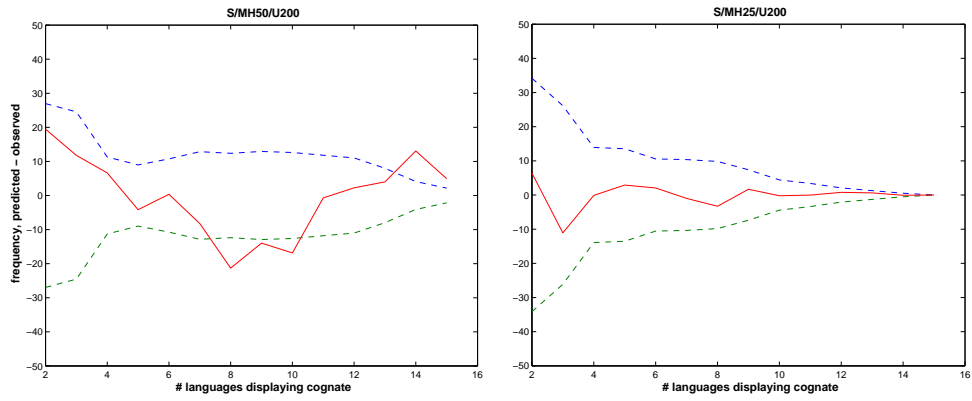


Fig. 15. Posterior predictive distributions (synthetic data) for $Y^{(n)}$, the number of traits displayed at exactly n leaves. y -axis $E(\tilde{Y}^{(n)}|D) - Y^{(n)}(D)$ and the envelope $\pm 2\text{std}(\tilde{Y}^{(n)}|D)$. x -axis $n = 2, 3, \dots, L$. (Left) Predictive distribution for $S/MH50/U200$ posterior. (Right) $S/MH25/U200$ posterior.

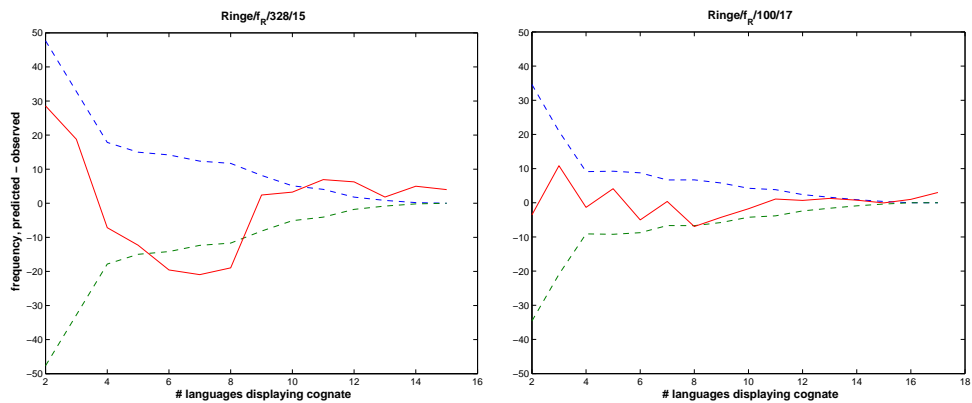


Fig. 16. Posterior predictive distributions (real data), otherwise as Fig. 15. (Left) Predictive distribution for $\text{Ringef}_R/328/15$ posterior. (Right) $\text{Ringef}_R/100/17$ posterior.