

relaxed phylogenetics and dating with confidence

Alexei Drummond

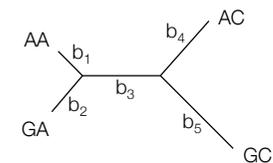
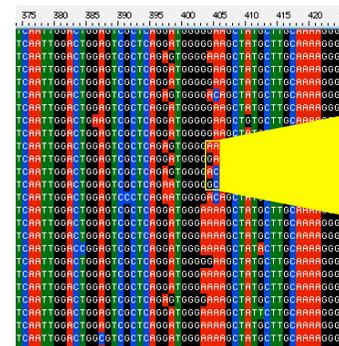
Overview

- Bayesian phylogenetics
- Relaxed molecular clocks in BEAST 1.4
- New relaxed molecular clock model

Review of Bayesian evolutionary inference

- The output of a Bayesian evolutionary analysis is a **probability distribution** on trees and parameter values.
- For **phylogenetics** the tree topology is the object of interest. The substitution parameters and tree prior parameters are a **nuisance** that we average over using MCMC and then ignore.
- For **population genetics** the tree and substitution parameters are a **nuisance** that we average over and then ignore, focusing instead on the population parameters.
- Sometimes an evolutionary hypothesis more specific than a full tree topology is of interest (like "Did this adaptive radiation predate the Miocene?") and then the result of the analysis should be the testing of this hypothesis, averaged over all trees and parameter values, weighted by their probability given the data.

Molecular evolutionary model: Felsenstein's likelihood (1981)

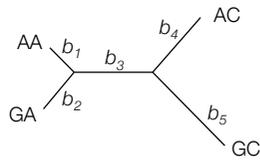


The probability of the sequence alignment,

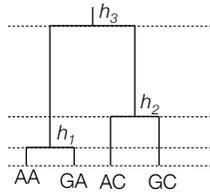
$$\Pr\{D | T, Q\}$$

can be efficiently calculated given a tree and branch lengths (T), and a probabilistic model of mutation represented by an instantaneous rate matrix (Q). In **phylogenetics**, branch lengths are usually unconstrained.

Model assumptions

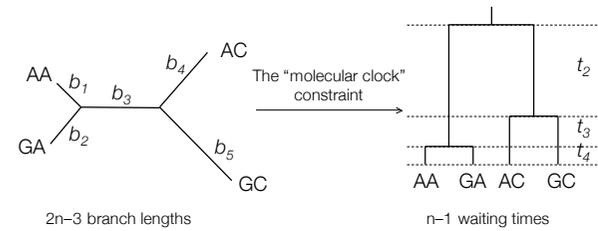


1. Product of rate and time (branch length) is independent and identically distributed among branches.
2. The root of the tree could be anywhere with equal probability.
3. Topology implies nothing about branch lengths.



1. Rate of evolution is the same on all branches.
2. The root of the tree is equidistant from all tips.
3. Topology constrains branch lengths (e.g. two branches in a cherry must be of equal length)

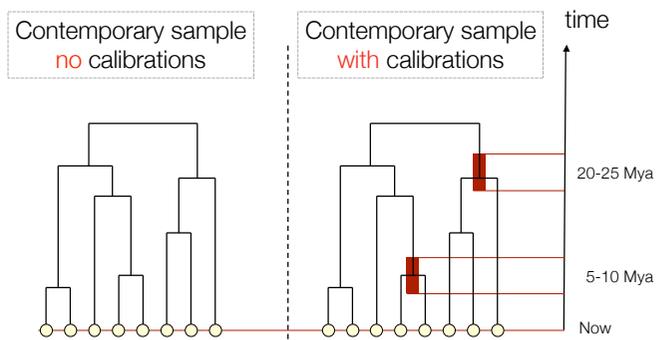
Combining the clock-constraint with Felsenstein's likelihood



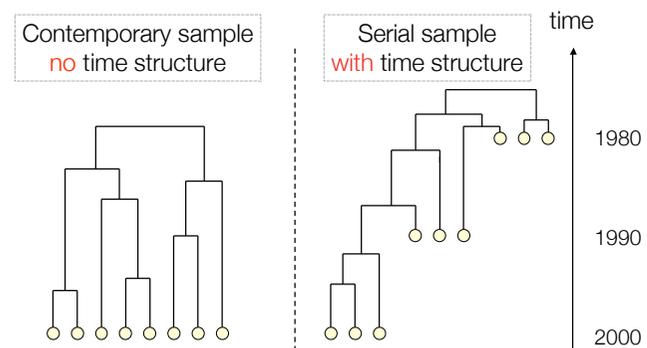
$$p(N, g, Q | D) \propto \Pr\{D | \mu, g, Q\} f_G(g | N) f_N(N) f_Q(Q)$$

The joint posterior probability of the tree (g), the tree prior parameters (N) and the substitution matrix (Q) are estimated using Markov chain Monte Carlo (Drummond et al, Genetics, 2002)

Time structure via calibrations



Time structure in samples themselves



Full Bayesian Model

Probability (density) of what we don't know given what we do know.

$$P(\mathbf{g}, \boldsymbol{\mu}, N_e, \mathbf{Q} | D) = \frac{1}{Z} P(D | \mathbf{g}, \boldsymbol{\mu}, \mathbf{Q}) f_G(\mathbf{g} | N_e) f_m(\boldsymbol{\mu}) f_N(N_e) f_Q(\mathbf{Q})$$

Likelihood function

other priors

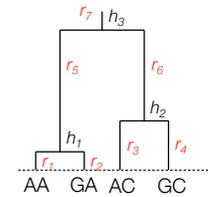
Unknown normalizing constant

tree prior

Q = substitution parameters
 N_e = tree prior parameters
 \mathbf{g} = tree
 $\boldsymbol{\mu}$ = overall substitution rate

In the software package BEAST, MCMC integration can be used to provide a chain of samples from this density.

Relaxing the molecular clock



In the field of **divergence time estimation** auto-correlated relaxed clocks have been considered.

e.g. Thorne et al, 1998:

$$r_i \sim \text{LogNormal}(r_{A(i)}, \sigma^2 \Delta t_i)$$

$$r \sim \text{Exp}(\lambda)$$

$$r \sim \text{LogNormal}(\mu, \sigma^2)$$

$$r \sim \text{Gamma}(\alpha, \beta)$$

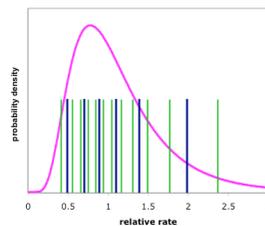
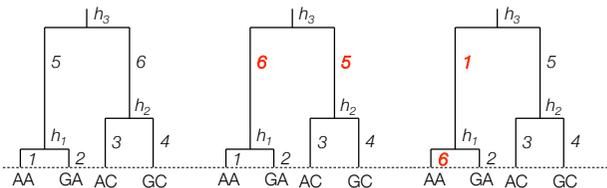
We introduce a relaxed clock model in which there is no prior correlation between child and parent rates

"Un-correlated" or "memory-less" relaxed clocks

ML

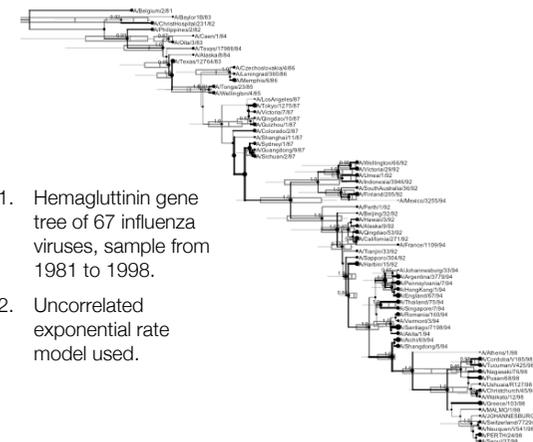
Drummond et al, available in BEAST

Sampling branch rates using MCMC

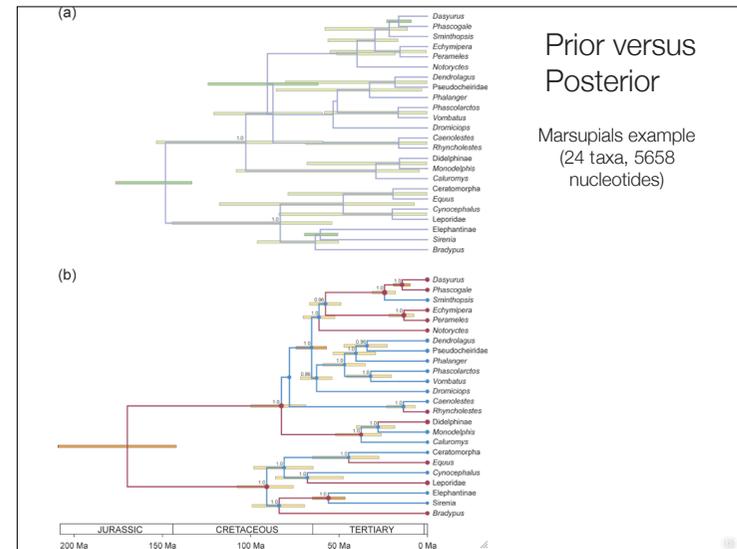
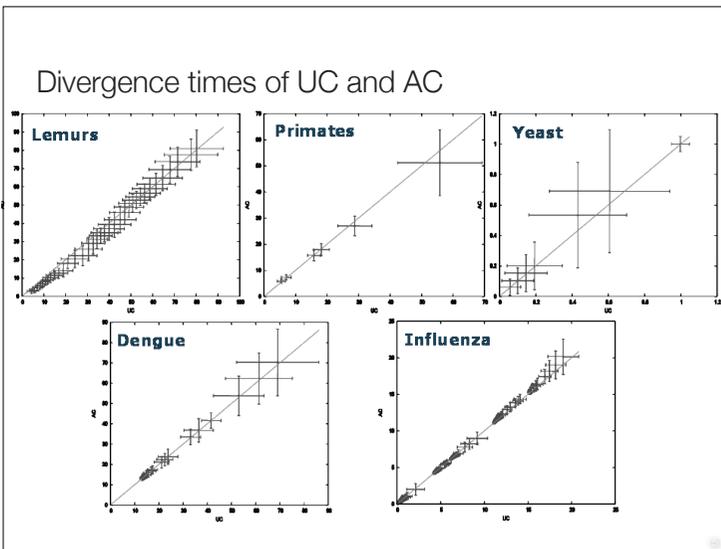
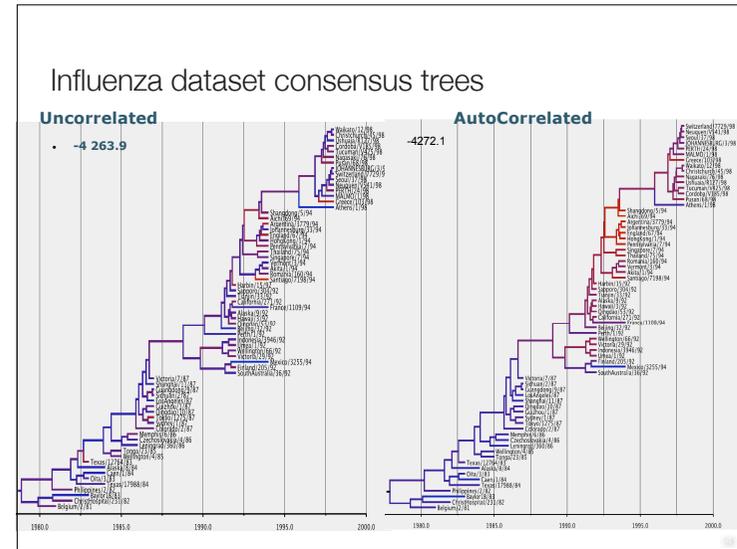
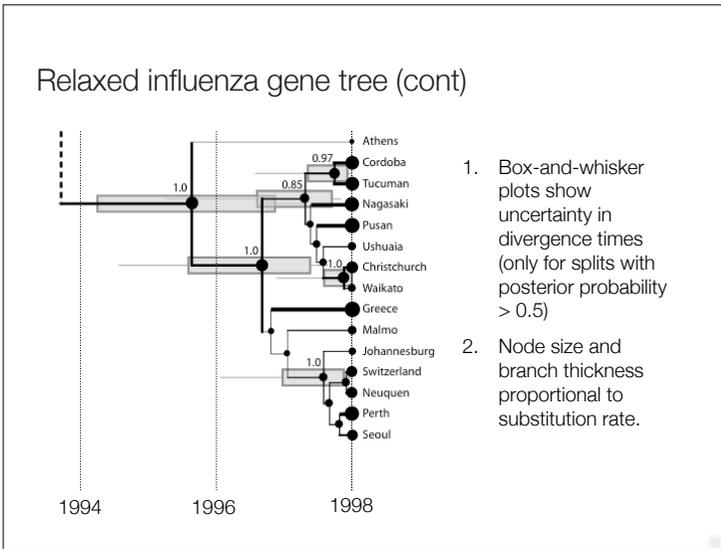


1. Rates are summarized into $2n-2$ rate categories (e.g. blue is 6 categories; green is 12 categories).
2. Random pairs of rates categories are swapped during MCMC.
3. For purposes of topology changes, rate categories are associated with child node.

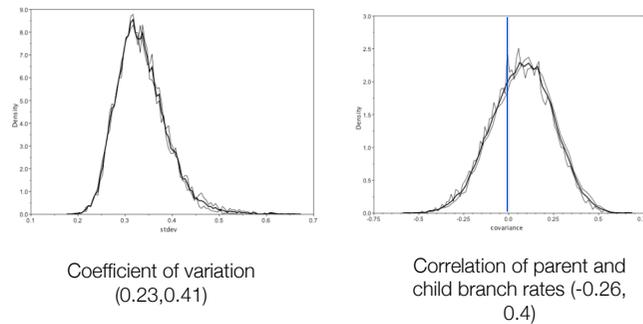
Relaxed influenza gene tree



1. Hemagglutinin gene tree of 67 influenza viruses, sample from 1981 to 1998.
2. Uncorrelated exponential rate model used.

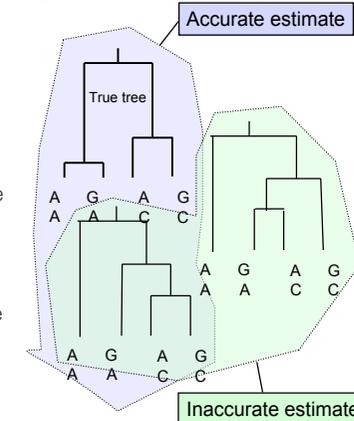


Posterior distribution of rates across branches



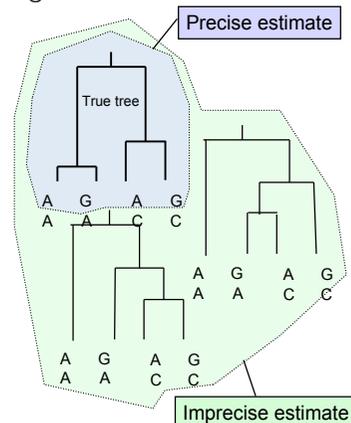
Accuracy in Bayesian phylogenetics

- Phylogenetics is an estimation problem, in which the phylogenetic tree topology is the object we wish to estimate.
- The error associated with this estimation can be described by the **95% credible set** of trees: the smallest set of trees including 95% of the posterior probability.
- A standard measure of accuracy is the **false positive rate**. How often do we exclude the true tree from the **95% credible set**? Ideally it would be 5%...



Precision in Bayesian phylogenetics

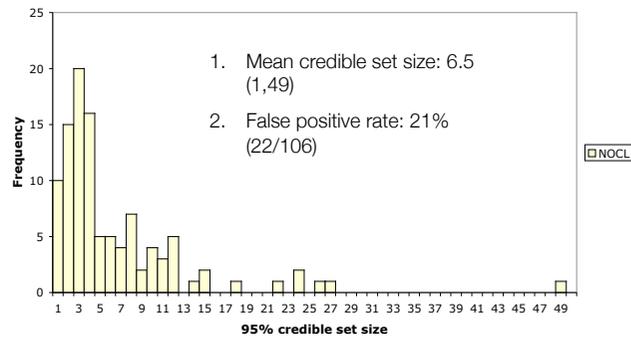
- The precision of an estimate can be described by how much is excluded.
- How small is the **95% credible set** of trees?



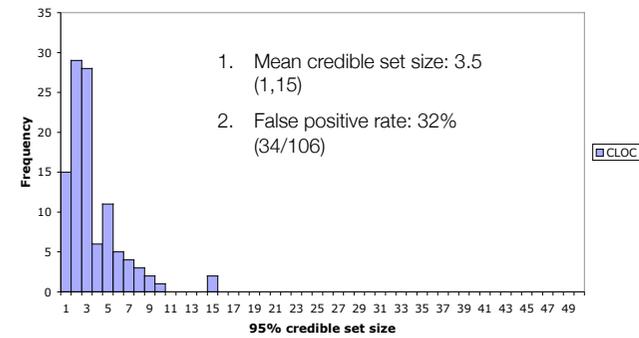
Testing accuracy and precision with real data

- Used 106 genes from 8 species of yeast (Rokas et al, 2003) and 4 other "phylogenomic" data sets
- For each gene used both MrBayes and BEAST to estimate phylogeny and 95% credible set
- Assumed true tree is the tree estimated using all the concatenated data set.
- Tabulated number of trees in credible set and whether the true tree was in credible set for MrBayes (unconstrained) and BEAST (MLLN and CLOC models)

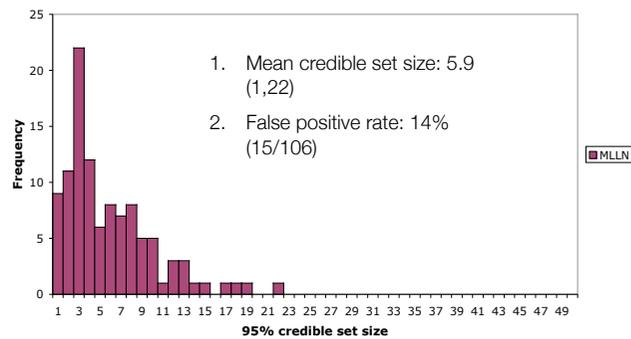
MrBayes results



Beast:CLOC results



Beast:UCLN results



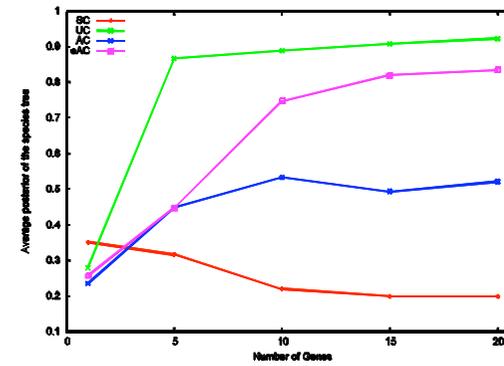
Summary of tree accuracy results

Dataset	Sample Size	Average Length	Clock Rejected by LRT	Accuracy (%) (True Tree in 95% Credible Set) ^a		
				CLOC	UCLN	UF
Bacteria	102	170 aa	26%	46.1	48.0	42.2
Yeast	106	1,198 bp	76%	67.0	84.9	79.2
Plants	61	647 bp	67%	91.8	88.5	83.6
Animals	99	197 aa	59%	64.6	69.7	57.6
Primates	500	632 bp	13%	88.8	89.0	88.8

Summary of tree precision results

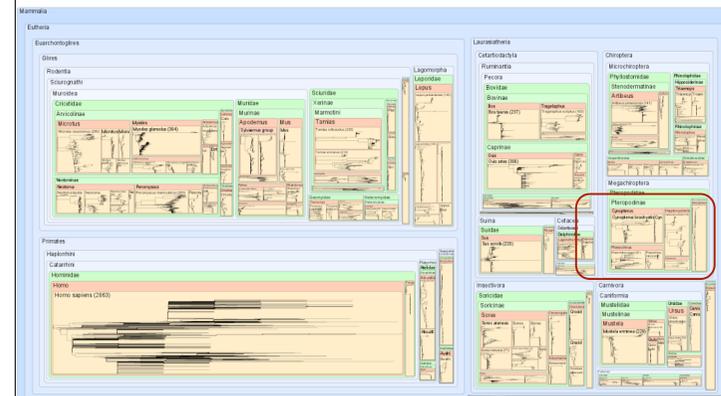
Dataset	Sample Size	Average Length	Clock Rejected by LRT	Precision (Number of Trees in 95% Credible Set) ^b		
				CLOC	UCLN	UF
Bacteria	102	170 aa	26%	5.7	10.3	11.3
Yeast	106	1,198 bp	76%	3.5	5.9	6.5
Plants	61	647 bp	67%	7.5	15.4	9.2
Animals	99	197 aa	59%	5.7	10.2	14.2
Primates	500	632 bp	13%	3.1	3.4	5.1

Increasing the length of the sequence



new relaxed clock model

Cytochrome b gene trees for >100 Mammals in Genbank

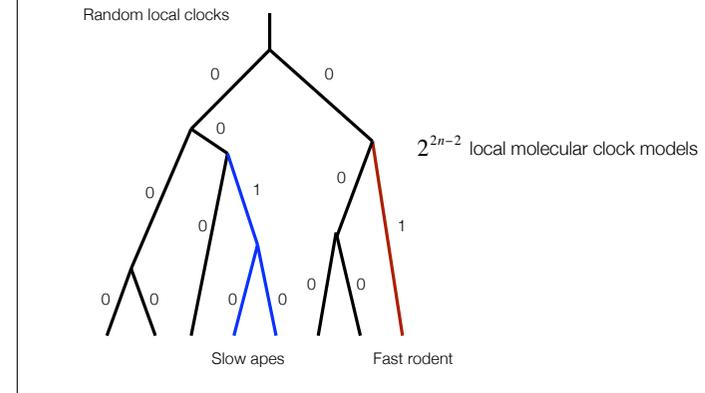


Check this BLAST view and more at <http://www.geneious.com>

When is different really different?



A new model for relaxing the clock



A new model for relaxing the clock

$$\mathbf{d} = \{\delta_1, \delta_2, \dots, \delta_{2n-2}\}$$

$$\delta_k \in \{0, 1\}$$

Indicators

$$\mathbf{f} = \{\phi_1, \phi_2, \dots, \phi_{2n-2}\}$$

$$\phi_i \in (0, \infty)$$

Rate scale parameters

$$\mathbf{r} = \{r_1, r_2, \dots, r_{2n-2}\}$$

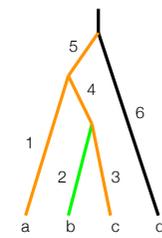
Rate scales

$$r_i = \phi_i^{\delta_i}$$

$$\mathbf{m} = \{\mu_1, \mu_2, \dots, \mu_{2n-2}\}$$

Branch rates

$$\mu_i = c(\mathbf{f})r_i^{\text{parent}(i)}$$



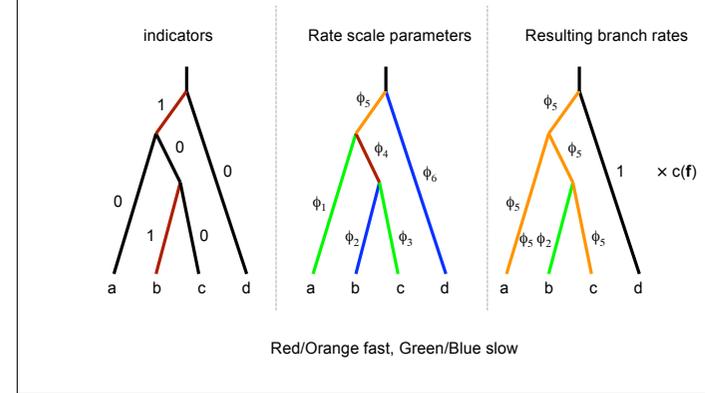
$$\mathbf{d} = \{0, 1, 0, 0, 1, 0\}$$

$$\mathbf{f} = \{\phi_1, \phi_2, \phi_3, \phi_4, \phi_5, \phi_6\}$$

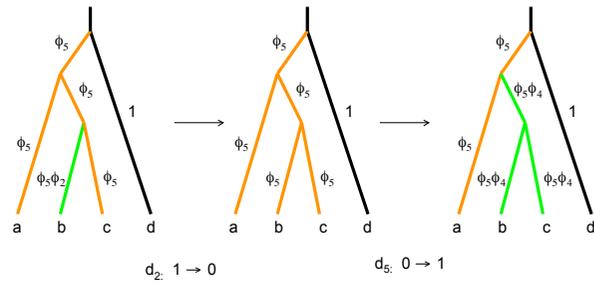
$$\mathbf{r} = \{1, \phi_2, 1, 1, \phi_5, 1\}$$

$$\mathbf{m} = c(\mathbf{f})\{\phi_5, \phi_2, \phi_5, \phi_5, \phi_5, 1\}$$

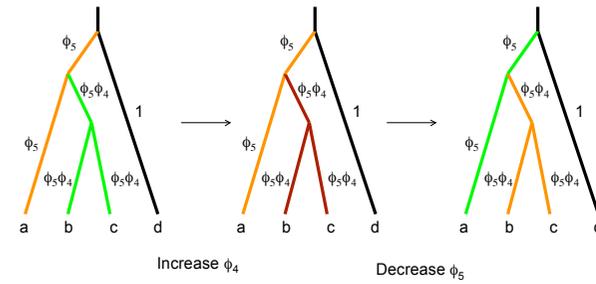
A new model of relaxed clock



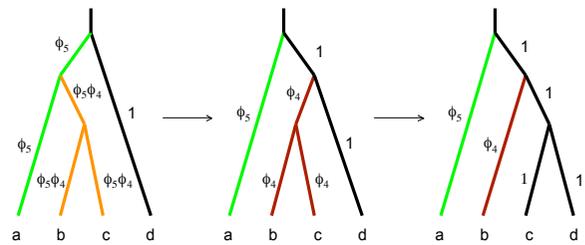
Sampling the indicators, d



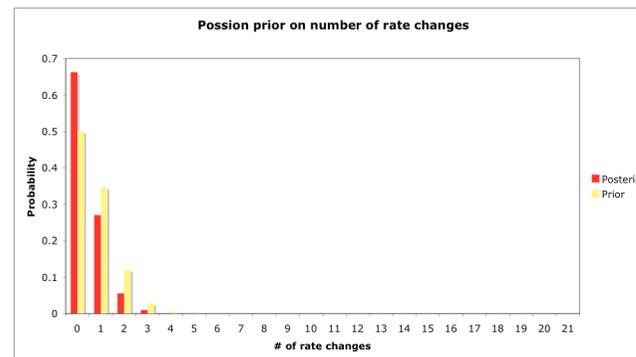
Sampling the rate parameters, f



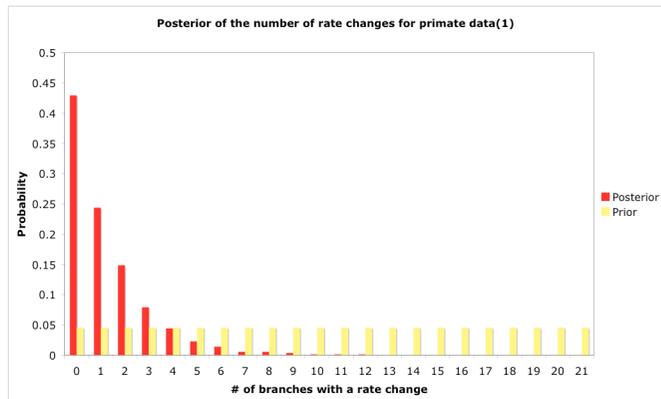
Sampling the trees



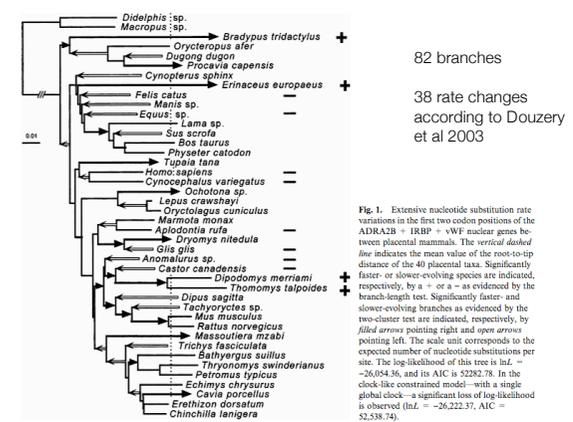
Primate.nex (Poisson prior)



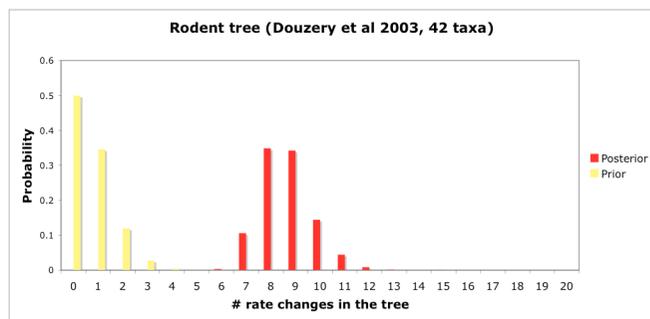
Primate.nex (Uniform prior)



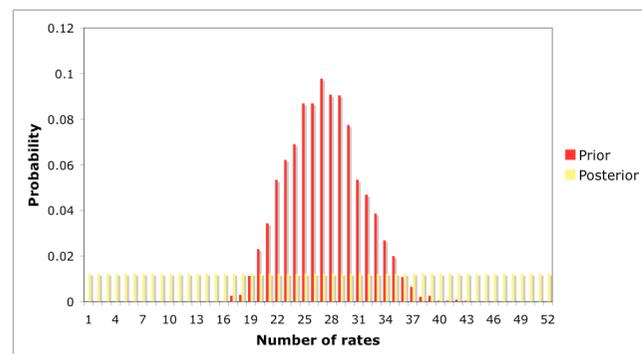
Rodents (1+2 pos of 3 nuclear genes)



Douzery et al, 2003 (Poisson prior)



Douzery et al, 2003 (Uniform prior)



Conclusions

- Relaxed molecular clocks have many benefits over unconstrained models for phylogenetic inference
 - › They appear to estimate the phylogenetic tree more accurately on real data sets
 - › They automatically provide estimates of a root, without the need for an outgroup
 - › They automatically provide estimates of relative divergence dates, or absolute divergence dates when calibration information is available
- Future directions
 - › Phylogenetic inference with autocorrelated rates (done, not published)
 - › Estimation of correlations in rate variation across multiple genes

the end

Phylogenetics: One tree to rule them all

- Majority consensus tree
(what branch lengths?)
Might not exist in the MCMC sample
- Highest posterior density state
tree/branches/parameters
Might just have very good branch lengths, but otherwise be a relatively unlikely topology
- Highest posterior density tree
(averaging branches/parameters)
Hard to estimate if there are many different trees in the 95% credible interval of trees
- Median tree?
Select the tree in the sample that minimizes the distance to the other trees using some metric.
- Maximum credibility tree?