

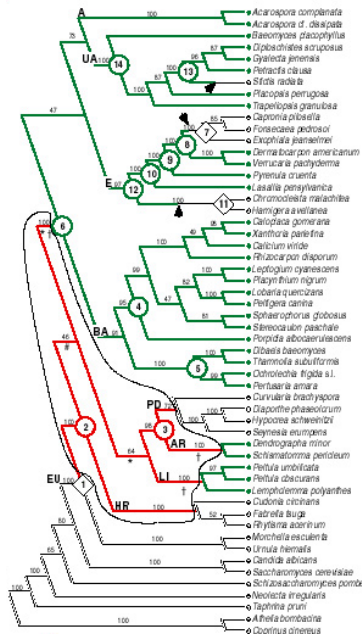
Accounting for Phylogenetic Uncertainty in Comparative Studies:
MCMC and MCMCMC Approaches

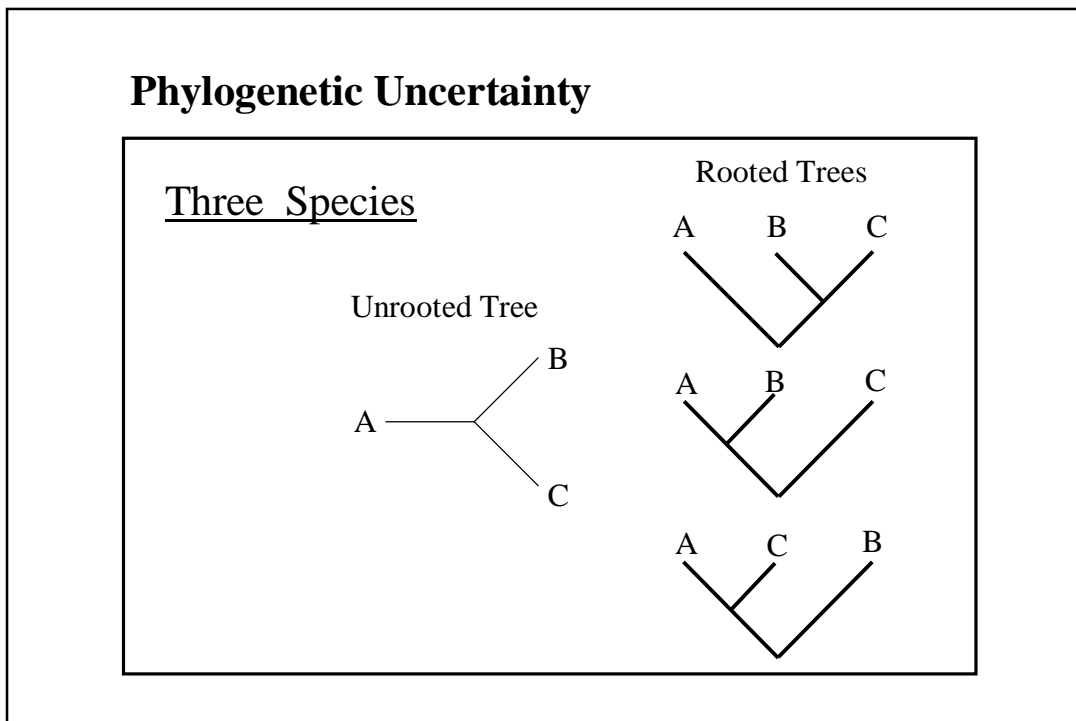
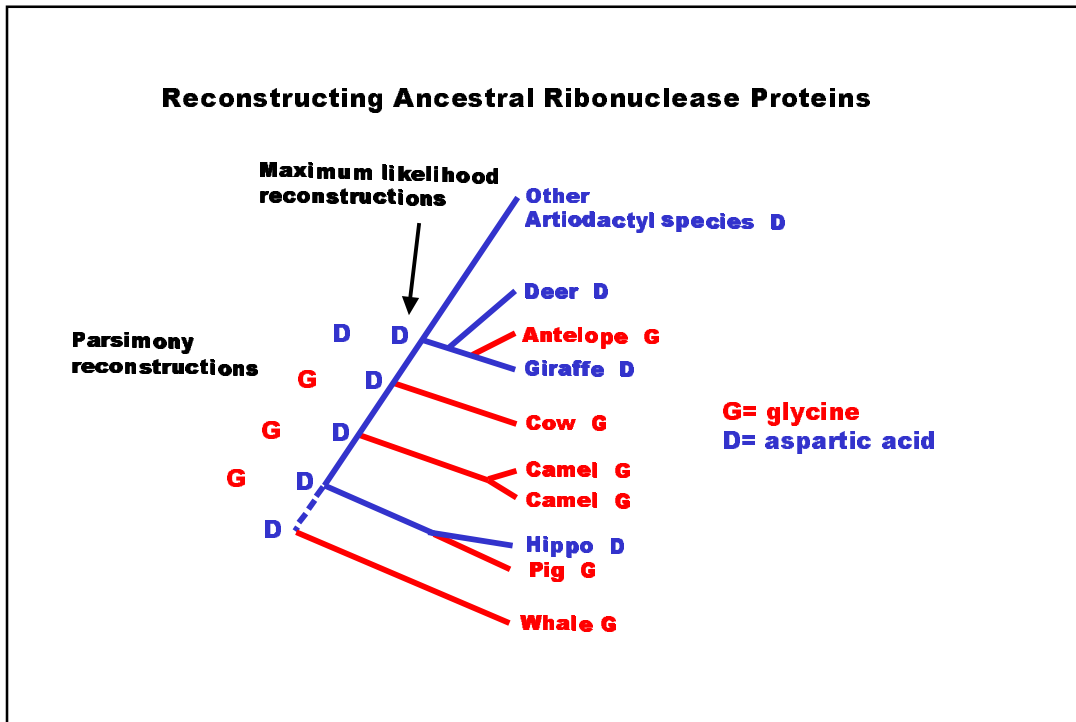
Mark Pagel
Reading University

m.pagel@rdg.ac.uk

Phylogeny of the Ascomycota Fungi
showing the evolution of
lichen-formation

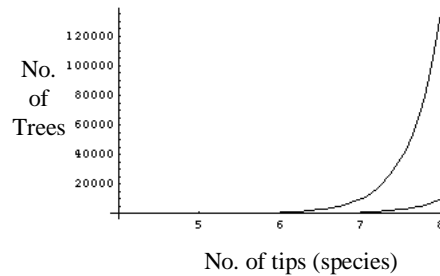
- █ lichen forming
- █ ambiguous
- █ not lichen forming





Number of Possible Phylogenetic Trees

Species	Unrooted	Rooted
3	1	3
4	3	15
5	15	105
6	105	945
10	2,027,025	34,459,425
50	2.83×10^{74}	2.75×10^{76}



N=50 No. rooted = 27529213532835651545259729751524430639300973035816196098326553772152587890625
 No. unrooted = 283806325080779912837729172696128150920628587998105114415737667754150390625

Accounting for Phylogenetic Uncertainty

Markov-Chain Monte Carlo (MCMC) Methods

- generate a long chain of phylogenetic trees (tree proposal and acceptance mechanisms)
- randomly sample from the converged chain
- calculate event or evolutionary process in each

Tree acceptance mechanism: The Metropolis-Hastings Algorithm

Accept new tree with $p=1.0$ if $L(T_{n+1}) > L(T_n)$

otherwise...

accept with probability $\propto L(T_{n+1}) / L(T_n)$

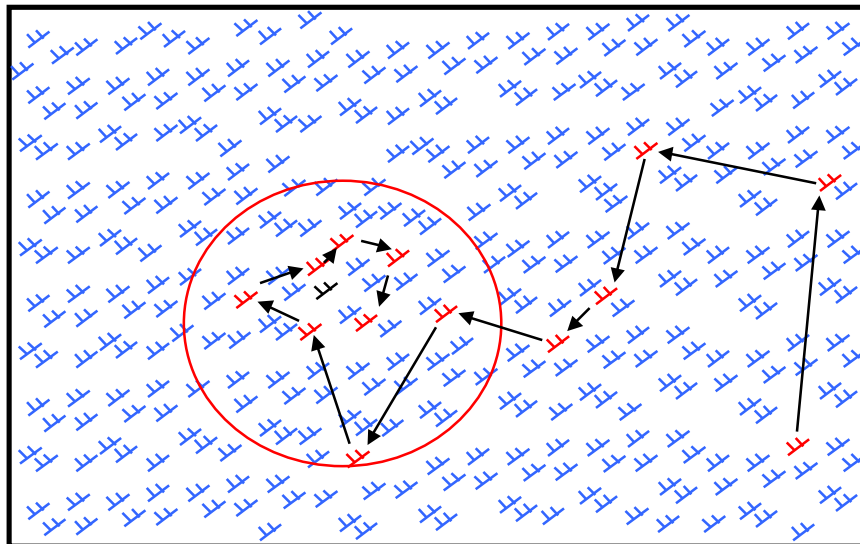
Metropolis-Hastings Algorithm: Accept new tree according to:

$$R = \min \left\{ 1, \frac{f(X|T')}{f(X|T)} \times \frac{f(T)}{f(T')} \times \frac{f(T|T')}{f(T|T)} \right\}$$

likelihood ratio prior ratio proposal ratio

X=data (e.g., gene sequences) T=tree (topology, branches, parameters)

MCMC Sampling



Primer of finding the likelihood of a phylogenetic tree

1. aligned gene- sequence data

Sheep	ATGGTGAAAA	GCCACATAGG	CAGTTGGATC	CTGGTTCTCT	TTGTGGCCAT
Human	ATGGCGAA--	----CCTTGG	CTGCTGGATG	CTGGTTCTCT	TTGTGGCCAC
Gorilla	ATGGCGAA--	----CCTTGG	CTGCTGGATG	CTGGTTCTCT	TTGTGGCCAC
Mink	ATGGTGAAAA	GCCACATAGG	CAGCTGGCTC	CTGGTTCTCT	TTGTGGCCAC

2. model of sequence evolution

3. the probability of sequence substitutions in a given branch of the phylogeny

$$P(t) = \text{Exp}[Qt]$$

4. the likelihood of a given phylogenetic tree

$$L = \prod_{\text{branches}} P(t) = \prod [\text{Exp}[Qt]]$$

5. search alternative topologies

Sheep	ATGGTGAAAA	GCCACATAGG	CAGTTGGATC	CTGGTTCTCT	TTGTGGCCAT
Human	ATGGCGAA--	----CCTTGG	CTGCTGGATG	CTGGTTCTCT	TTGTGGCCAC
Gorilla	ATGGCGAA--	----CCTTGG	CTGCTGGATG	CTGGTTCTCT	TTGTGGCCAC
Mink	ATGGTGAAAA	GCCACATAGG	CAGCTGGCTC	CTGGTTCTCT	TTGTGGCCAC

$$p_i(x | T, v, Q) = \sum_n^{4^{s-1}} \left[w_{root(i)} \left(\prod_{k=1}^s p_{n_k, x_{ki}}(v_k, Q) \right) \left(\prod_{k=1}^{s-2} p_{n'_k, n_k}(v_k, Q) \right) \right]$$

↑

probability of observing
ith nucleotide

↑

prior weight
ith site

↑

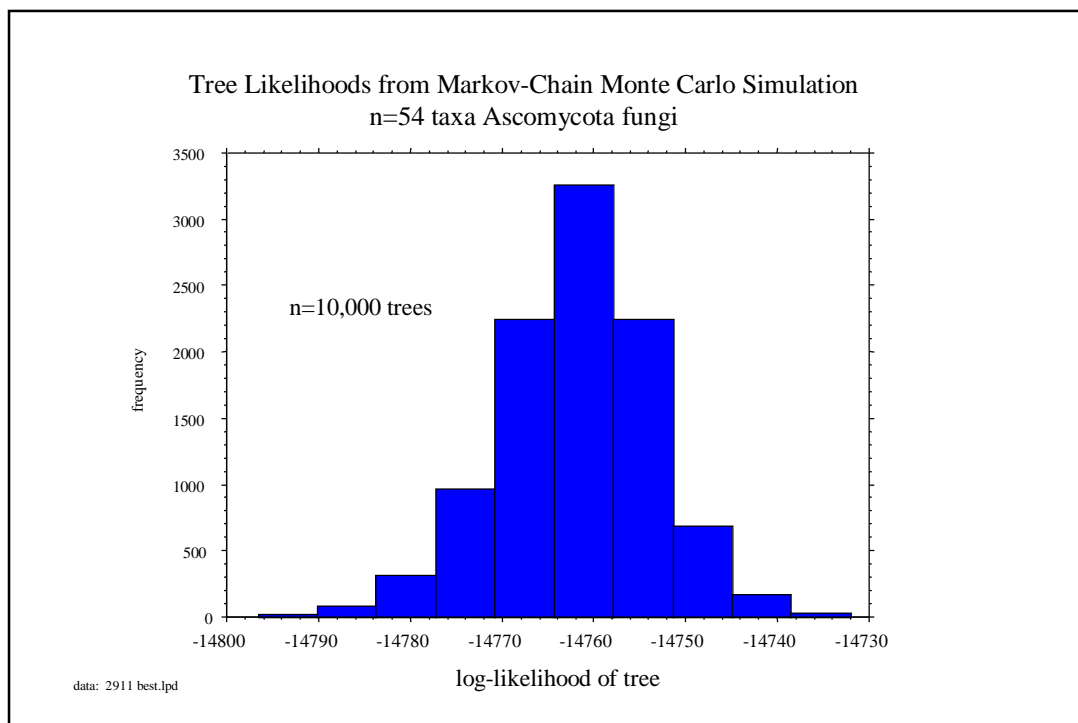
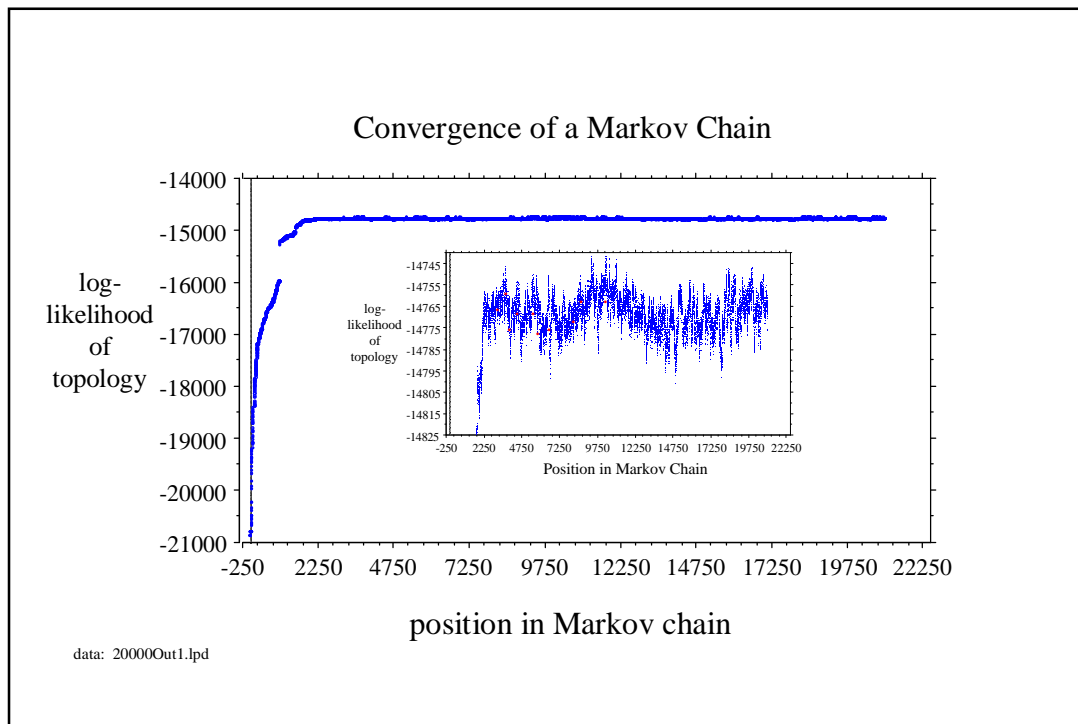
product over s branches
leading to species

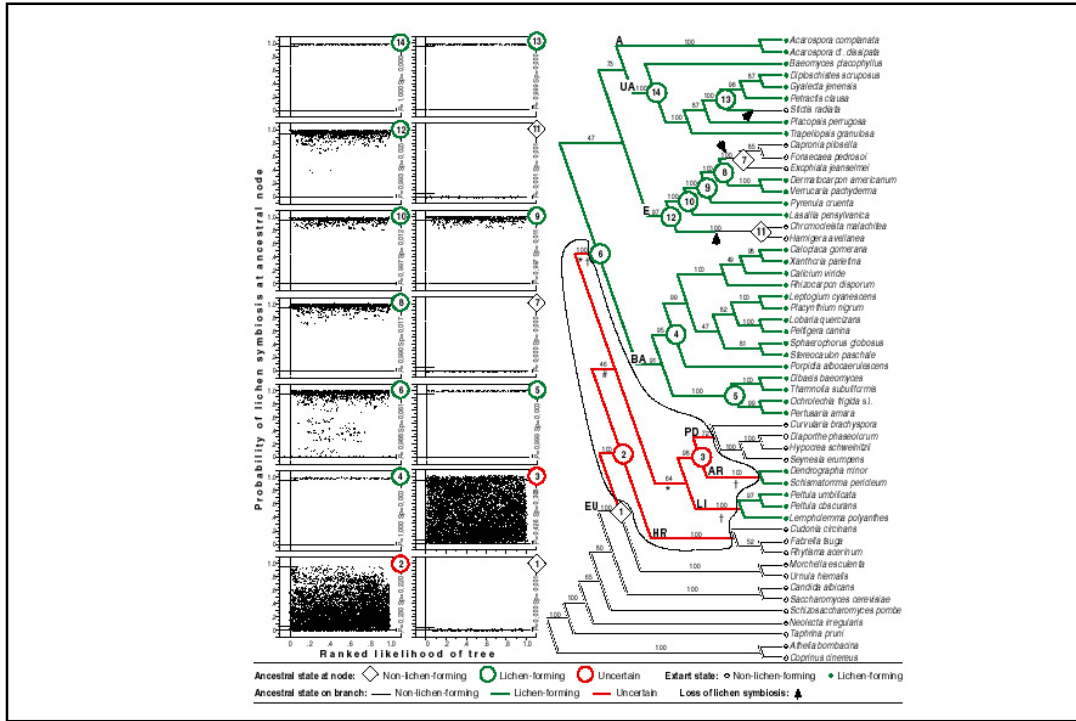
↑

product over s-2 internal
branches

possible assignments of
ancestral nodes (e.g., 64)

$$L(x | T, v, Q) = \prod_i p_i(x | T, v, Q) \leftarrow \text{product over all } i \text{ nucleotides}$$

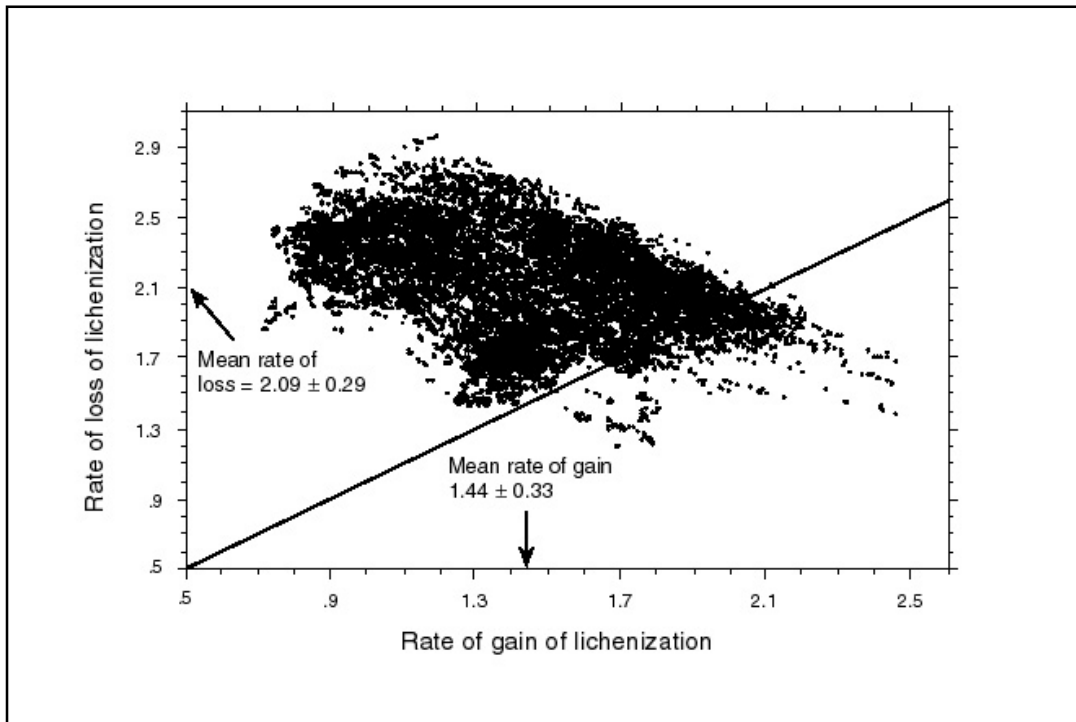




Character Transition Rates “gains” and “losses” of lichenization

	Single Tree*	MCMC ‘Integration’**
gains (q_{01}) =	1.04 (0.05-4.5)	1.47±0.32
losses (q_{10}) =	2.41(0.7-5.6)	2.12 ±0.31

*consensus tree **n=20,000 trees



MCMC Some Issues

- Lack of convergence (poor 'mixing')
- Tree and parameter proposal mechanisms
- Tree and parameter updating schedules
- Detecting convergence

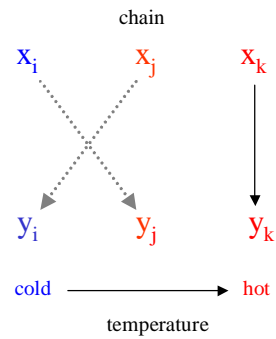
Metropolis-Coupled Markov Chain Monte Carlo (MCMCMC)

Given m simultaneous Markov chains, update chains, then swap states among a randomly chosen pair i and j each iteration according to:

$$R = \min \left\{ 1, \frac{f_i(y_i) f_j(y_j)}{f_i(x_i) f_j(x_j)} \right\}$$

{likelihood ratio chain i * likelihood ratio chain j }

probability of swapping with chain $i = R * 2/m$



'Temperatures' of heated chains

