

Statistical frameworks for modelling in phylogenetics

Maximum likelihood

Maximum Likelihood Estimation (MLE)



Joe Felsenstein

- Statistical method for estimating parameters of a model (e.g. mean and variance of a normal distribution)
- Originally developed by R. A. Fisher in the 1920s
- Adapted for phylogenetics by Joe Felsenstein

Born	Joseph Felsenstein May 9, 1942 (age 81)
Alma mater	University of Chicago
Known for	PHYLIP Felsenstein's tree-pruning algorithm

[Wikipedia](#)

EVOLUTIONARY TREES FROM DNA-SEQUENCES - A MAXIMUM-LIKELIHOOD APPROACH

FELSENSTEIN, J

1981 | [JOURNAL OF MOLECULAR EVOLUTION](#) 17 (6) , pp.368-376

10,183

Citations

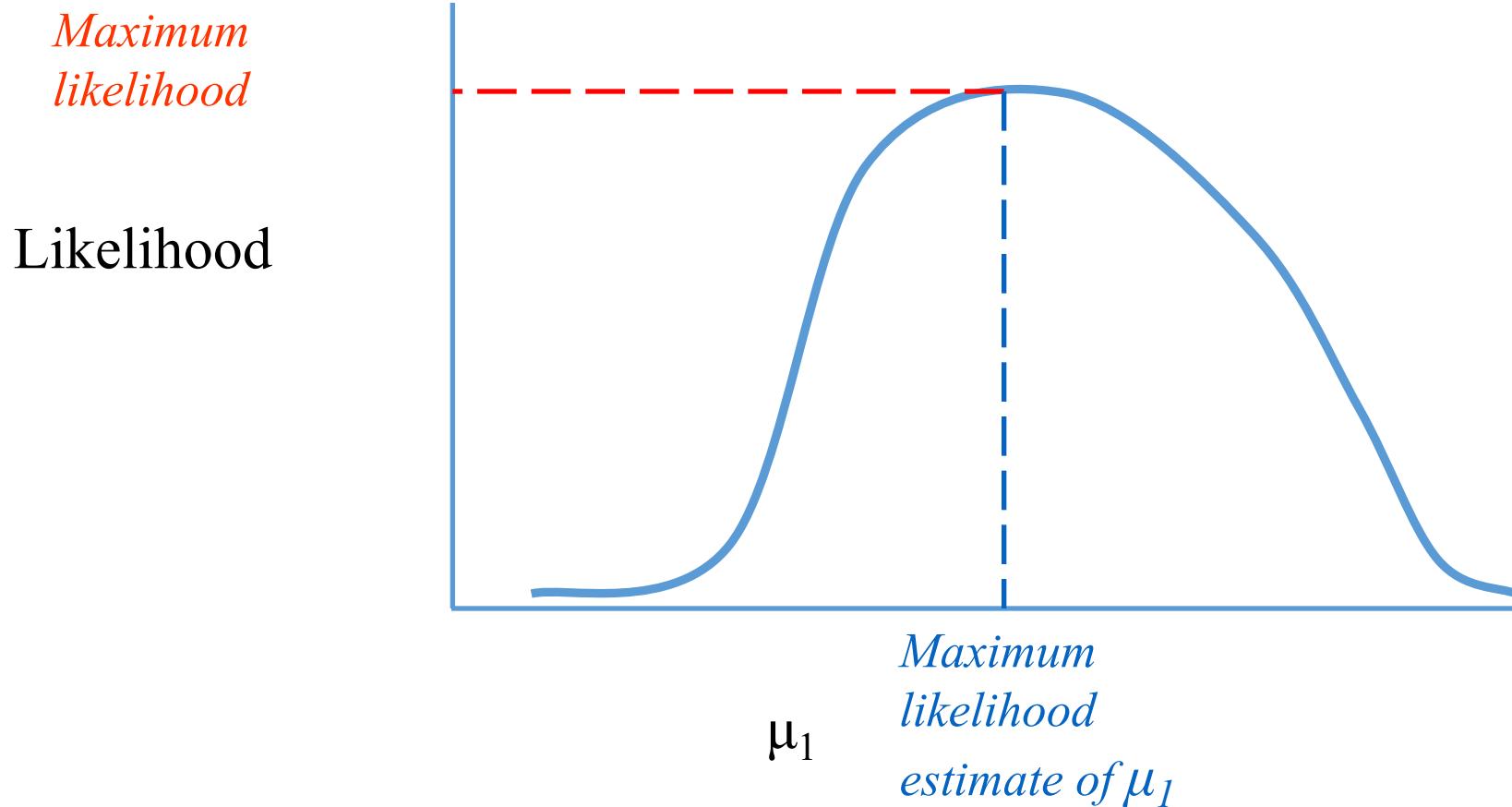
26

References

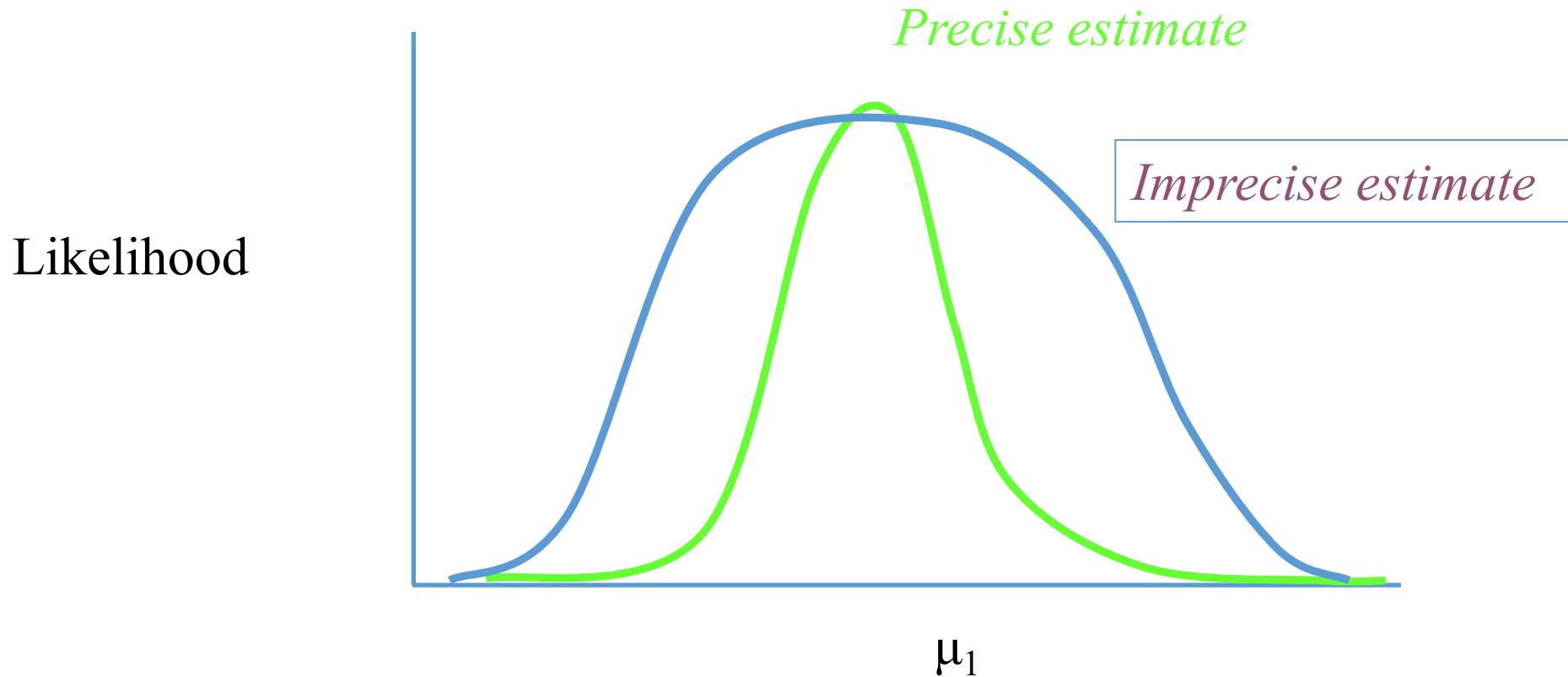
Likelihood of a hypothesis

- Likelihood (L) is proportional to the probability (Pr) of observing the data (D) given a model (M) – *conditional probability*
 - $L(M) = \text{Pr}(D | M)$
- We can examine this likelihood function to find where it is highest and identify the parameters of the model at this point -> Maximum Likelihood Estimates

Maximum Likelihood



Maximum Likelihood



Likelihood of a hypothesis

- Likelihood (L) is proportional to the probability (P) of observing the data (D) given a model (M) – *conditional probability*
 - $L(M) = \Pr(D | M)$
- We can examine this likelihood function to find where it is highest and identify the parameters of the model at this point -> Maximum Likelihood Estimates
- In molecular phylogenetics, likelihood is the **probability of observing the sequences given our model** (e.g. GTR+G and our tree topology)

Maximum Likelihood

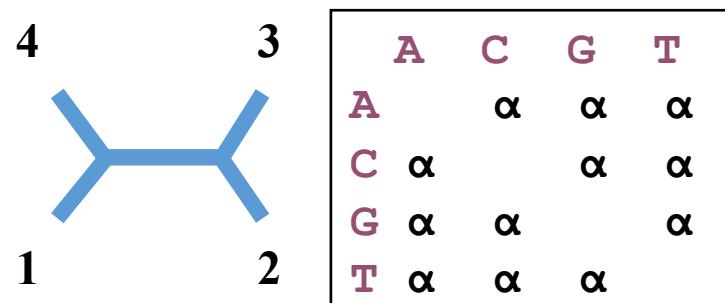
- For reconstructing phylogenies

Model

- which tree topology (τ), branch lengths, and parameters of DNA evolution model (θ) (e.g. transition/transversion ratio, base frequencies, ...) are maximizing the probability of observing the sequences at hand?

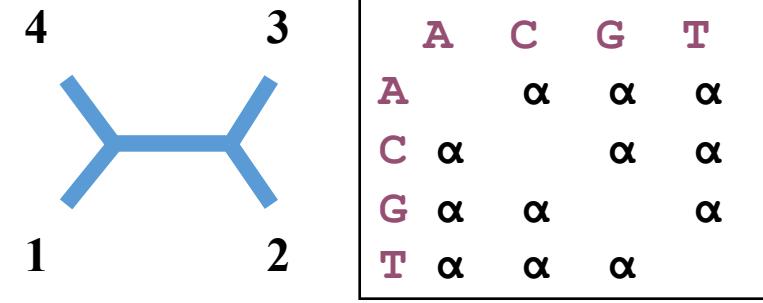
Data

$$L(\tau, \theta) = \Pr(\text{Data} \mid \tau, \theta)$$



=

AAGTTTGTGATTGCTCCCGTCATTA
AAGTTTGTGATTATTACCGGCCATTCATTA
AAGATTCTGATTATTACCCCATTCATTA
AAGTTCTGATTATTACCTCCATTCTTTA
AAGTTTTGATTACTCCCCGGTCTCTTA
AAGATTTCGGTTACTACCCCATCACTA
AAGATTTCGATTATTGCCCTTCATTA
AAGATTTTGATTATTACCTCCATTCTTTA

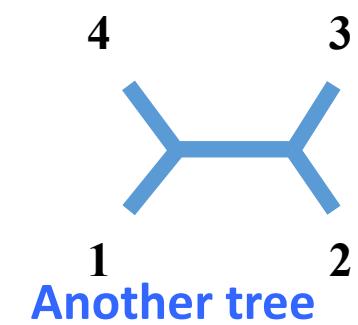
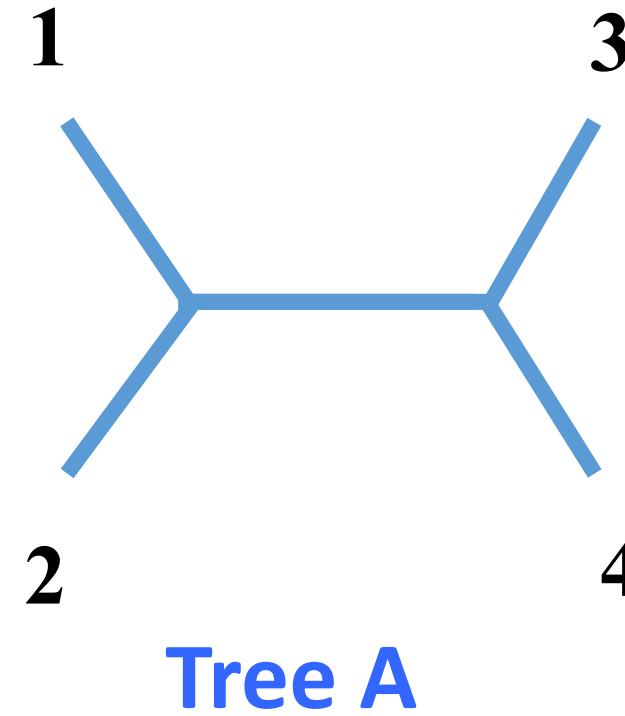


ML analysis in short

- Tree topology is obtained
- Branch lengths and parameters of the DNA substitution model are optimized
- Different topologies (with branch lengths and DNA substitution model parameters optimized) are compared based on their likelihood as the optimality criterion
- The topology with the highest likelihood needs to be found

Maximum Likelihood tree reconstruction

1 CGAGAA
2 AGCGAA
3 AGATTC
4 GGATAT

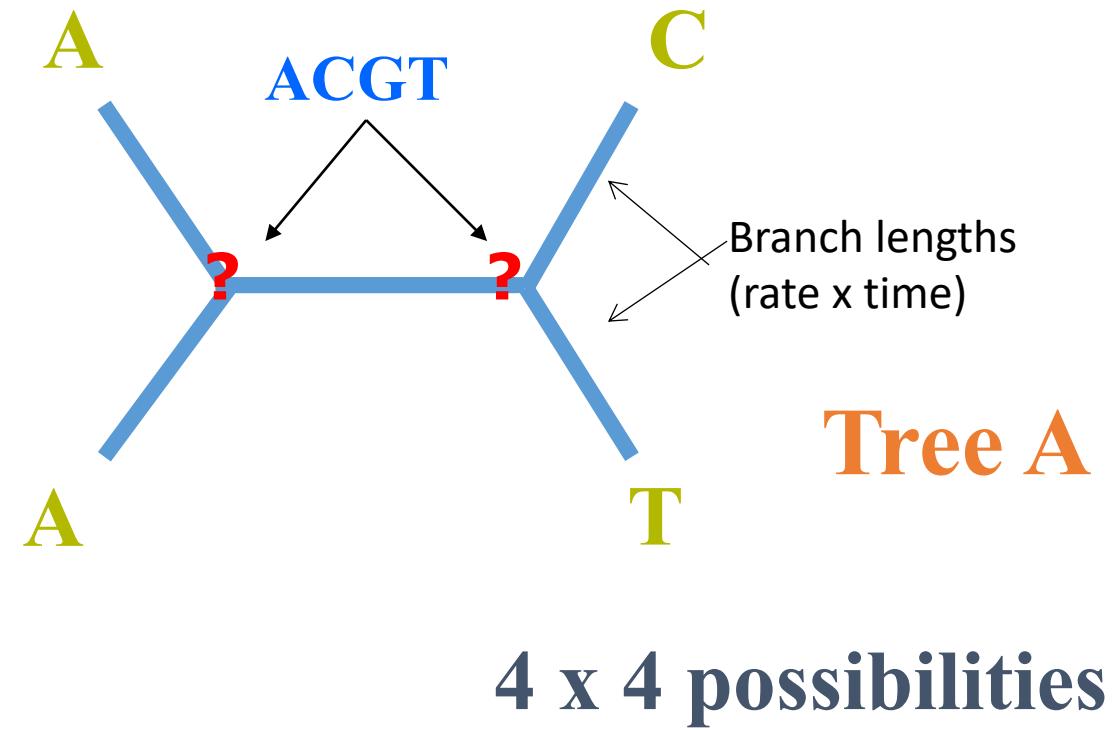


What is the likelihood that Tree A (rather than another tree) could have generated the sequence alignment?

Maximum Likelihood tree reconstruction

1	CGAGA	A
2	AGCGA	A
3	AGATT	C
4	GGATA	T

j

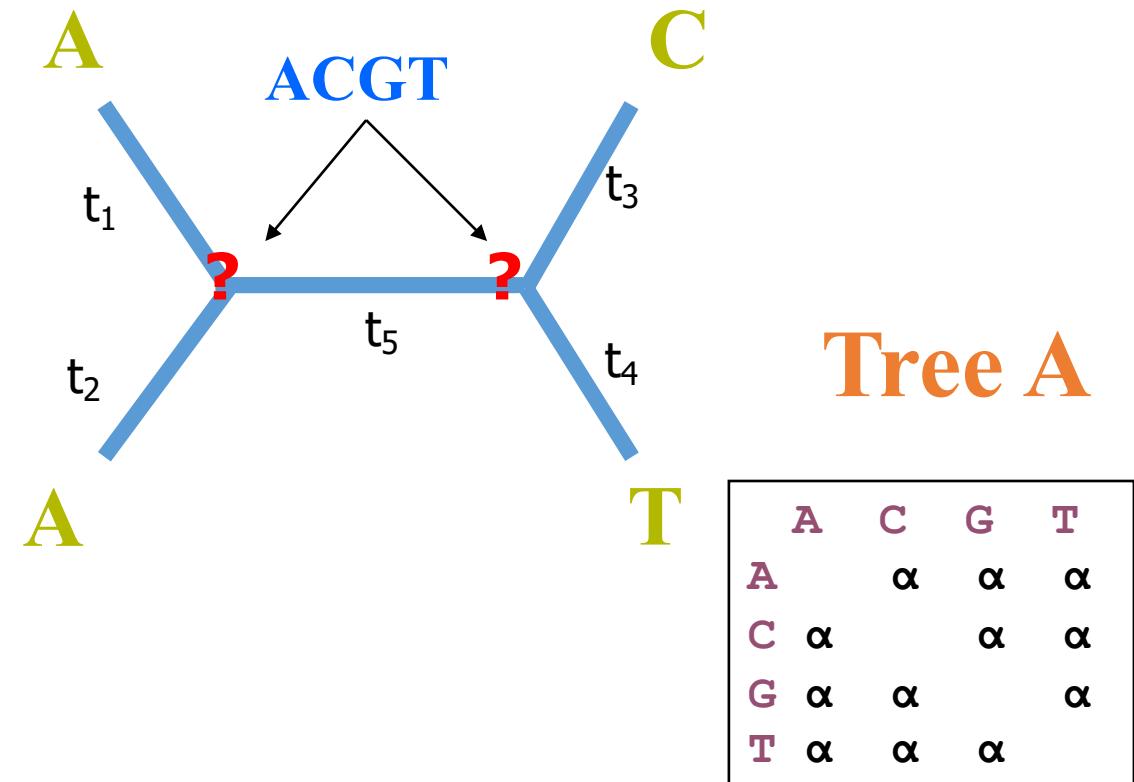


The likelihood for a particular site j is the sum of the probabilities of every possible reconstruction of ancestral states under a chosen model

Maximum Likelihood tree reconstruction

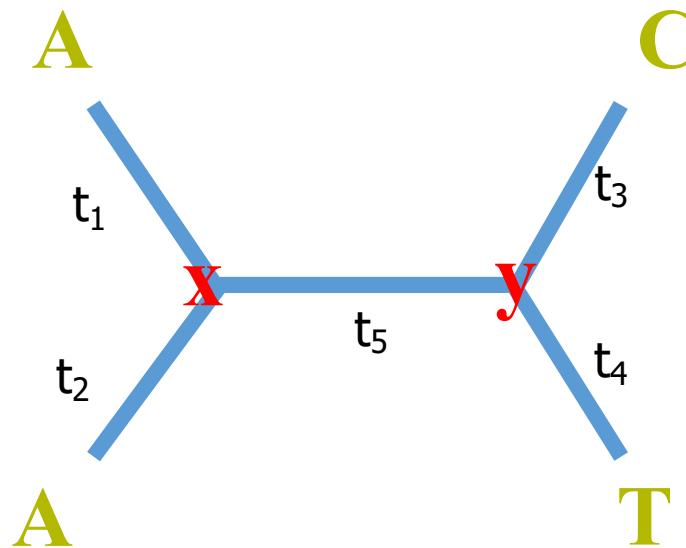
1	CGAGA	A
2	AGCGA	A
3	AGATT	C
4	GGATA	T

j



The likelihood for a particular site j is the sum of the probabilities of every possible reconstruction of ancestral states under a chosen model

Branch lengths also need to be estimated!



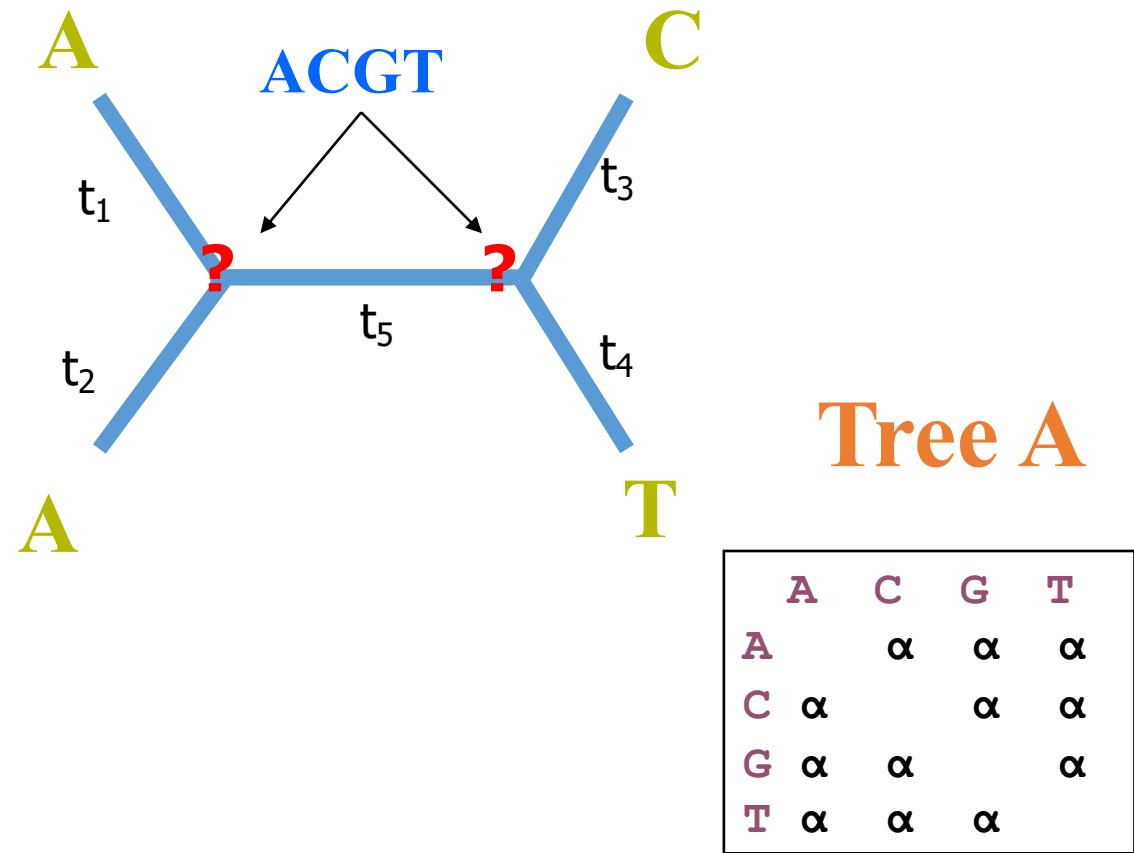
t_i are branch
lengths
(rate \times time)

$$P(A, A, C, T, x, y | T) = \text{Prob}(x) \text{ Prob}(A|x, t_1) \text{ Prob}(A|x, t_2) \text{ Prob}(y|x, t_5) \text{ Prob}(C|y, t_3) \text{ Prob}(T|y, t_4)$$

Maximum Likelihood tree reconstruction

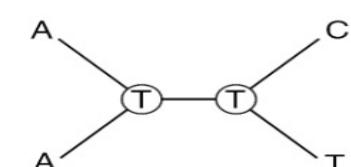
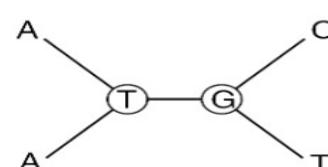
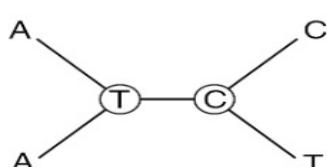
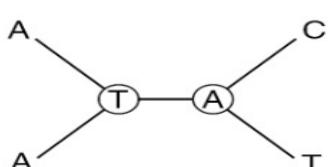
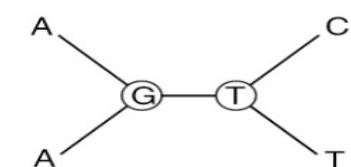
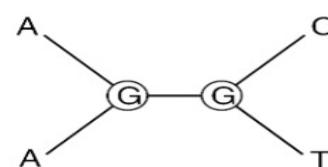
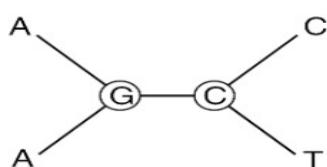
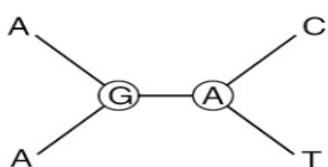
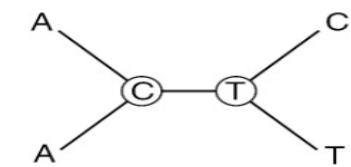
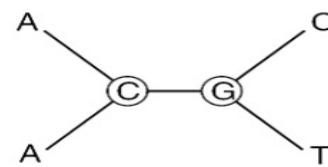
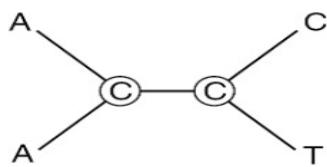
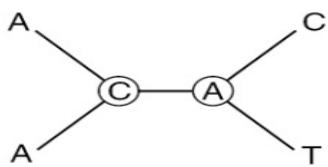
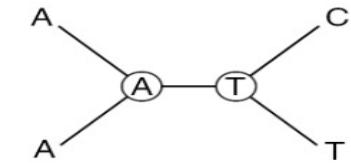
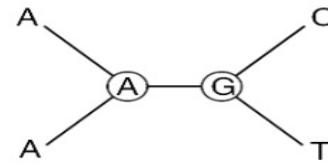
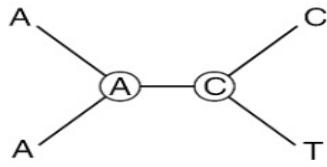
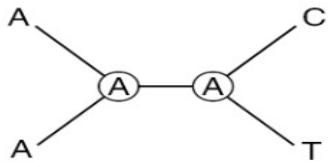
1	CGAGA	A
2	AGCGA	A
3	AGATT	C
4	GGATA	T

j



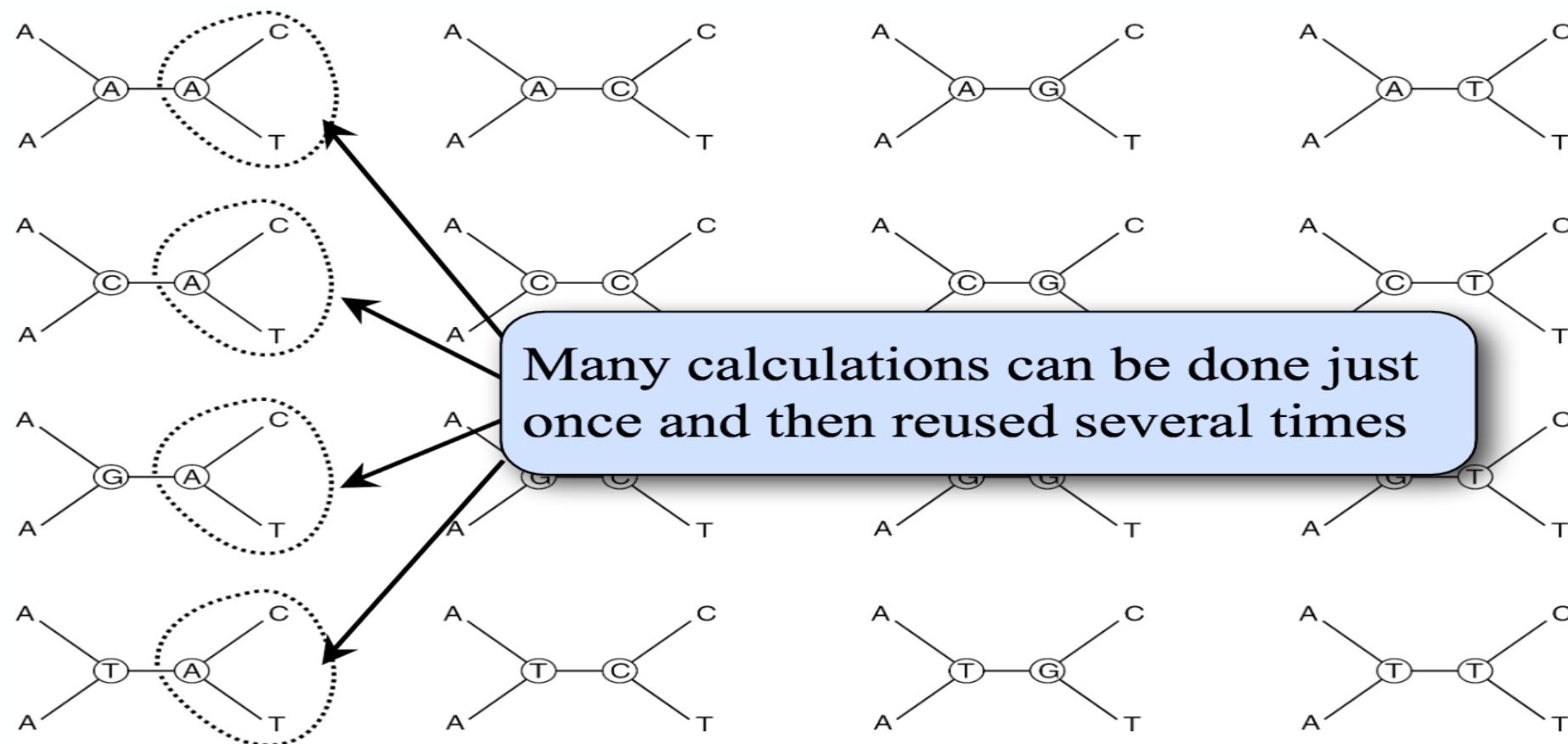
The likelihood for a particular site j is the sum of the probabilities of every possible reconstruction of ancestral states under a chosen model

Brute force approach would be to calculate L_k for all 16 combinations of ancestral states and sum them



Note use of the OR probability rule

Pruning algorithm (same result, less time)

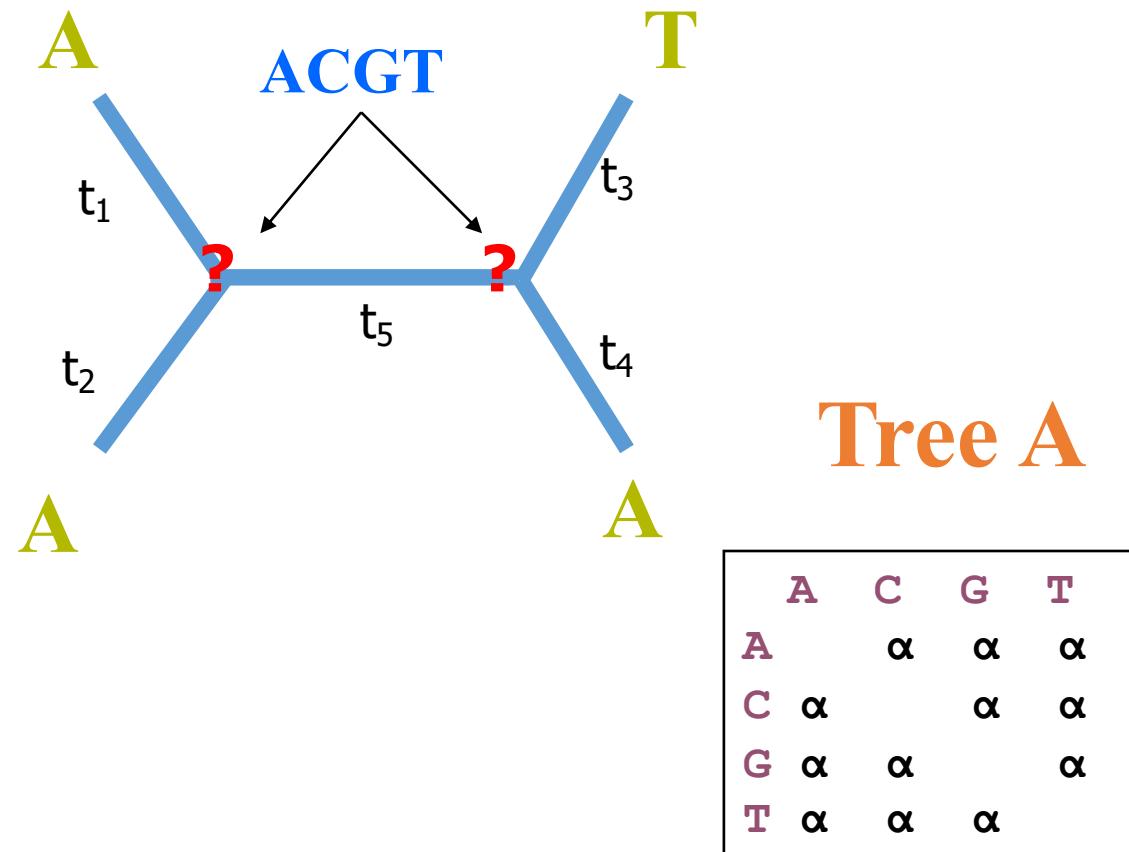


Felsenstein, J. 1981. Evolutionary trees from DNA sequences:
a maximum likelihood approach. *Journal of Molecular Evolution* 17:368-376

Maximum Likelihood tree reconstruction

1	CGAG	A	A
2	AGCG	A	A
3	AGAT	T	C
4	GGAT	A	T

j



The likelihood for a particular site j is the sum of the probabilities of every possible reconstruction of ancestral states under a chosen model

ML analysis in summary

- The likelihood for each site in the alignment for a given topology, branch lengths and the DNA substitution model is calculated
- The probability of the single site is the sum of probabilities of each scenario, taking into account all of the possible nucleotides that may have existed as states at the internal nodes
- The likelihood for a given tree topology for the whole alignment is the product of the likelihoods for each site

Typical assumptions of ML substitution models

- The probability of any change is independent of the prior history of the site (**a Markov Model**)
- Relative frequencies of A, G, C, and T are at equilibrium (**stationarity - S**)
- Change is **time reversible** - R e.g. the rate of change of A to T is the same as T to A
- Substitution probabilities do not change with time or over the tree (**a homogeneous Markov process – H**)
- **SRH** – we assume that sequence evolution is stationary, reversible, and homogeneous or SRH

Finding the maximum likelihood of a tree

- Problem: the number of possible trees (e.g. for 10 taxa, 2 million unrooted trees possible; for 60 taxa, more possible trees than atoms in the universe!)
 - for each tree topology we need to identify the maximum likelihood estimate for evolutionary parameters and branch lengths
 - then compare the likelihood among all the trees
 - This is simply computationally not feasible
- ▶ Solution:
 - Currently no method guarantees finding the best tree
 - Starting tree made usually using MP or NJ
 - Heuristic approaches are used:
 - e.g. NNI = Nearest Neighbour Interchange, SPR = subtree pruning and regrafting, TBR = tree-bisection and reconnection

Numbers of possible trees for N taxa

1	1
2	1
3	1
4	3
5	15
6	105
7	945
8	10395
9	135135
10	2027025
11	34459425
12	654729075
13	13749310575
14	316234143225
15	7905853580625
16	213458046676875
17	6190283353629370
18	191898783962510625
19	6332659870762850625
20	221643095476699771875 (2×10^{20})
50	3×10^{74}

Finding the maximum likelihood of a tree

- Problem: the number of possible trees (e.g. for 10 taxa, 2 million unrooted trees possible; for 60 taxa, more possible trees than atoms in the universe!)
 - for each tree topology we need to identify the maximum likelihood estimate for evolutionary parameters and branch lengths
 - then compare the likelihood among all the trees
 - This is simply computationally not feasible
- ▶ Solution:
 - Currently no method guarantees finding the best tree
 - Starting tree made usually using MP or NJ
 - Heuristic approaches are used:
 - e.g. NNI = Nearest Neighbour Interchange, SPR = subtree pruning and regrafting, TBR = tree-bisection and reconnection

A Bayesian Approach to Phylogenetics

A Bayesian approach compared to ML

- The likelihood is the probability of observing the data given a hypothesis
 - $L = \Pr(D | \theta)$.
- In ML we search for the parameter values of the model that maximize the likelihood function
- In a Bayesian analysis, we get the probability of a hypothesis given the data (probability of the tree given the sequences)
 - We combine the likelihood of a given hypothesis with a prior expectation for this hypothesis to obtain a posterior probability of the hypothesis

Bayes' rule in statistics

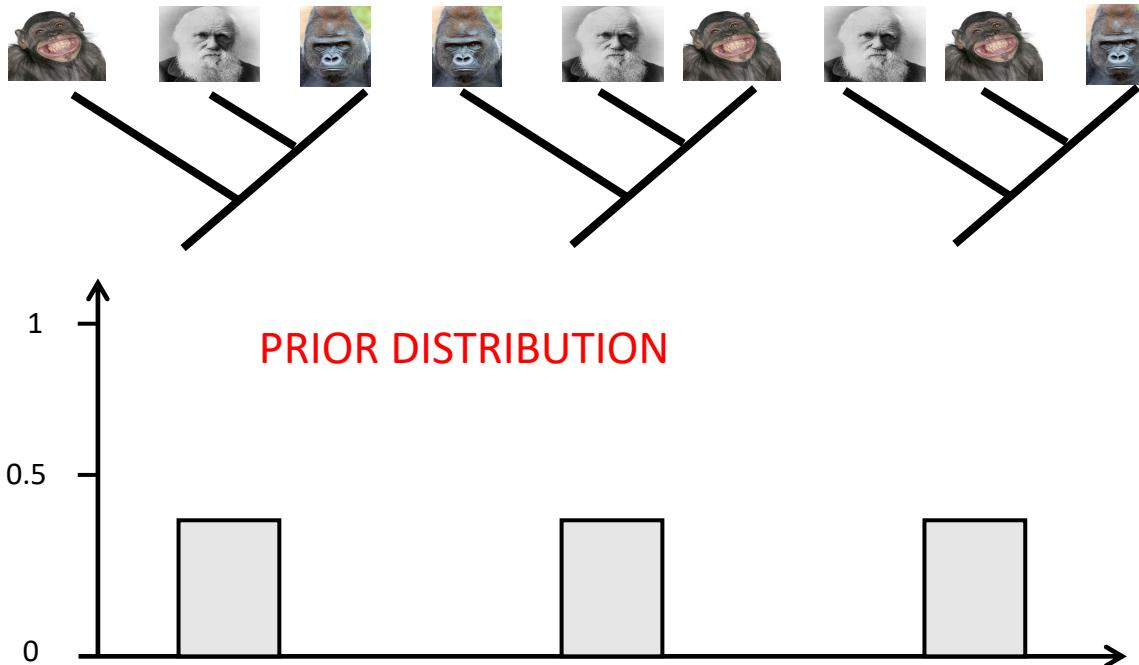
$$\Pr(\theta|D) = \frac{\Pr(D|\theta) \Pr(\theta)}{\sum_{\theta} \Pr(D|\theta) \Pr(\theta)}$$

Diagram illustrating Bayes' rule components:

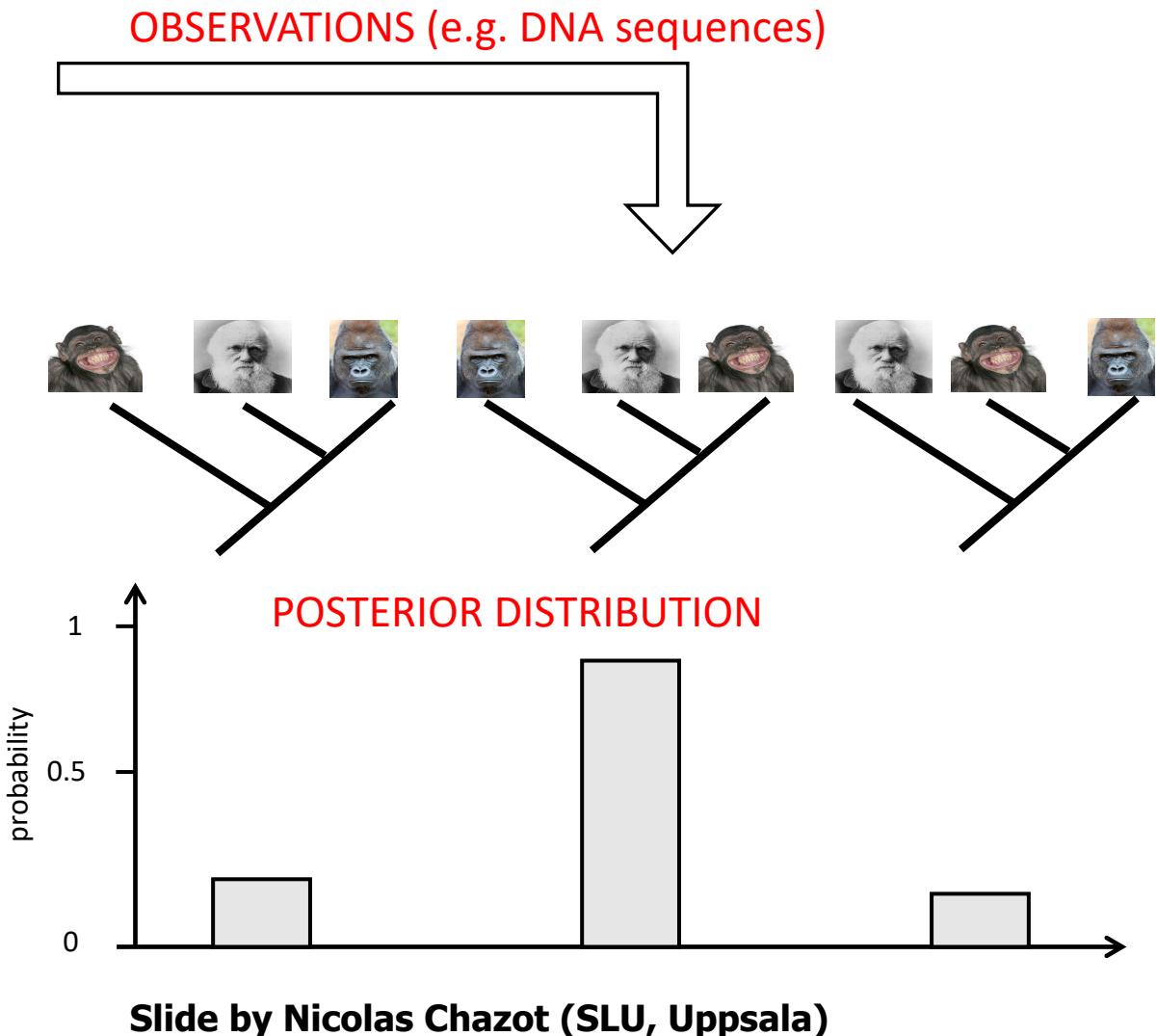
- Likelihood of hypothesis θ** : Points to the term $\Pr(D|\theta)$ in the numerator.
- Prior probability of hypothesis θ** : Points to the term $\Pr(\theta)$ in the numerator.
- Posterior probability of hypothesis θ** : Points to the result $\Pr(\theta|D)$.
- Marginal probability of the data (marginalizing over hypotheses)**: Points to the denominator $\sum_{\theta} \Pr(D|\theta) \Pr(\theta)$.

Bayesian inference in general

- D stands for data
- Θ (Gr. theta) means any one of a number of things:
 - a discrete hypothesis
 - a distinct model (e.g. JC, HKY, GTR, etc.)
 - a tree topology
 - one of an infinite number of continuous model parameter values (e.g. ts:tv rate ratio)
- Prior vs. posterior probability
- Posterior probability can be calculated by multiplying the prior probability of a tree (and model parameters) and the likelihood of the observed data (given a tree and its parameters) divided by a normalizing constant



- So-called flat priors for tree topology are used
 - All topologies have the same prior probability
 - Differences in posterior probability result from differences in the likelihood



Major difference between ML and BI

- In **ML joint estimation of parameters** – likelihood for all parameters estimated at once
 - Likelihood of each parameter depends on likelihood estimation of every other parameter
- In **BI marginal estimation** – posterior probability of any one parameter is calculated independently
- So even if using flat priors and the same model of DNA evolution, **ML and BI could infer different trees** because of differences between joint and marginal likelihood estimation

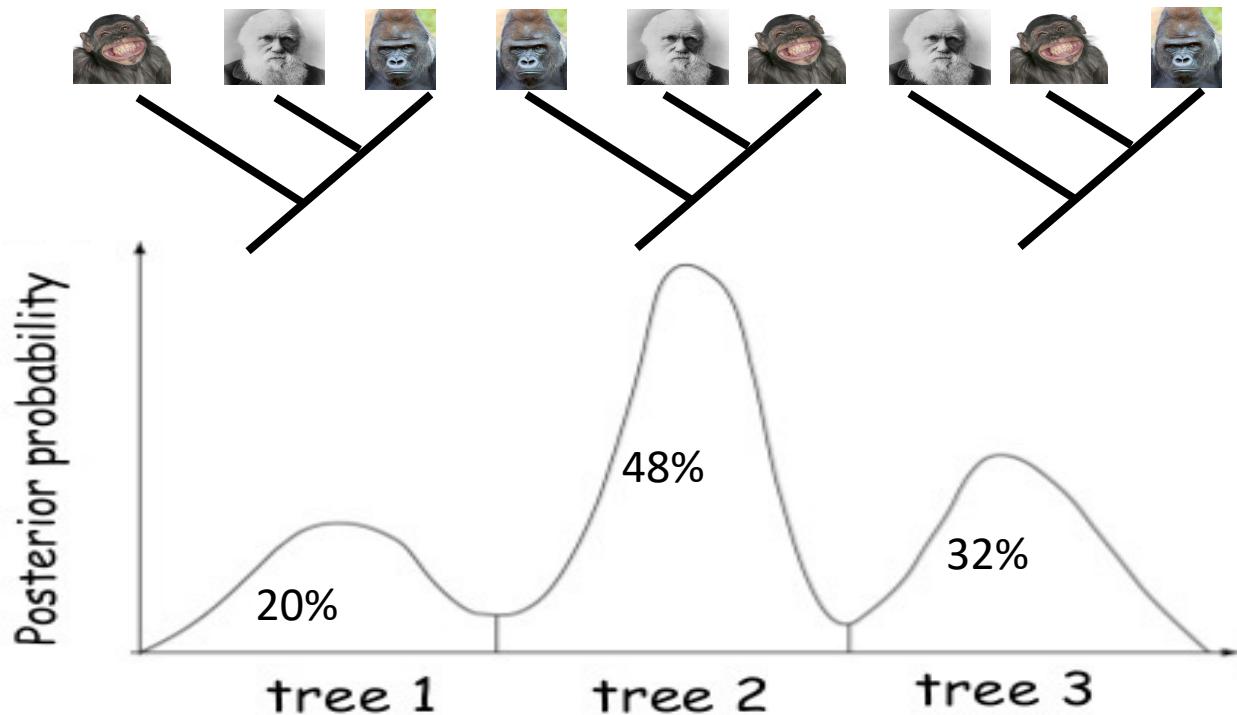
Bayes' rule: continuous case

$$f(\theta|D) = \frac{f(D|\theta)f(\theta)}{\int f(D|\theta)f(\theta)d\theta}$$

Diagram illustrating the components of Bayes' rule:

- Likelihood: $f(D|\theta)$ (blue box)
- Prior probability **density**: $f(\theta)$ (orange box)
- Posterior probability **density**: $f(\theta|D)$ (purple box)
- Marginal probability of the data: $\int f(D|\theta)f(\theta)d\theta$ (green box)

Arrows point from the Likelihood and Prior probability density boxes to the numerator of the equation. Arrows point from the Posterior probability density and Marginal probability of the data boxes to the terms in the denominator.



		Topologies			
		t1	t2	t3	
Branch length	vA	0.10	0.07	0.12	0.29
	vB	0.05	0.22	0.06	0.33
vC	0.05	0.19	0.14	0.38	
	0.20	0.48	0.32		

Joint probabilities

Marginal probabilities

Problem: it is impossible, in most cases,
to derive the posterior probability analytically

or even estimate it by drawing random samples from it

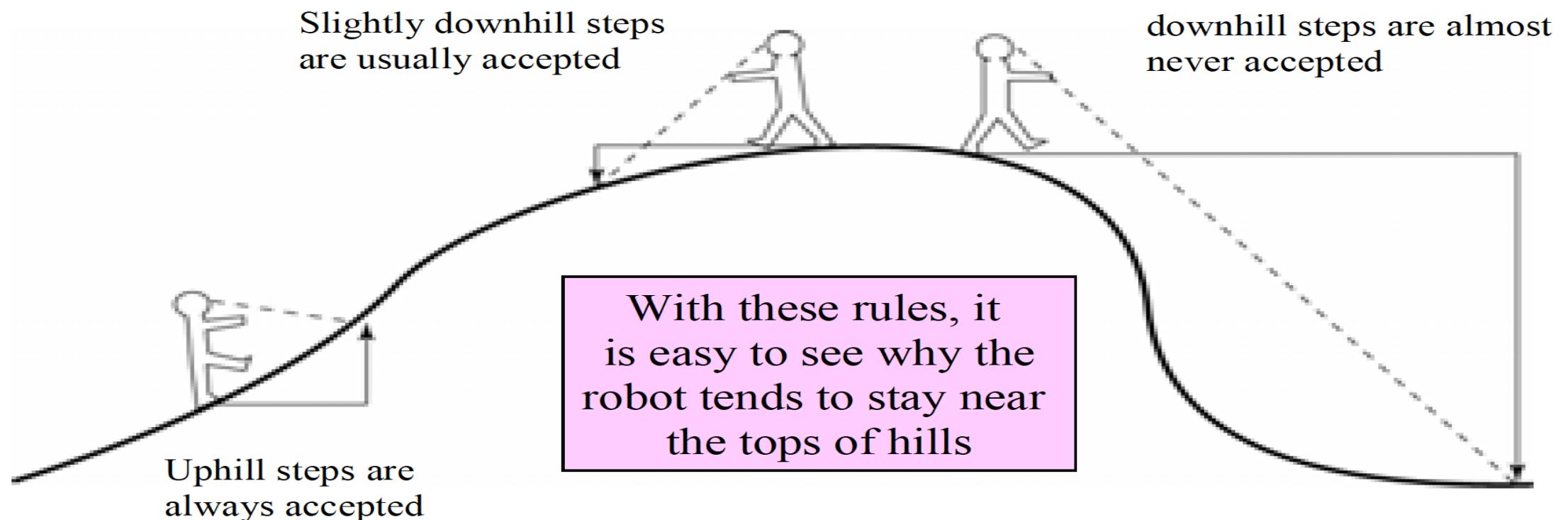
We want something that will “walk” across this parameter space and actively search for the highest point in the parameter “landscape”

Markov chain Monte Carlo (MCMC)

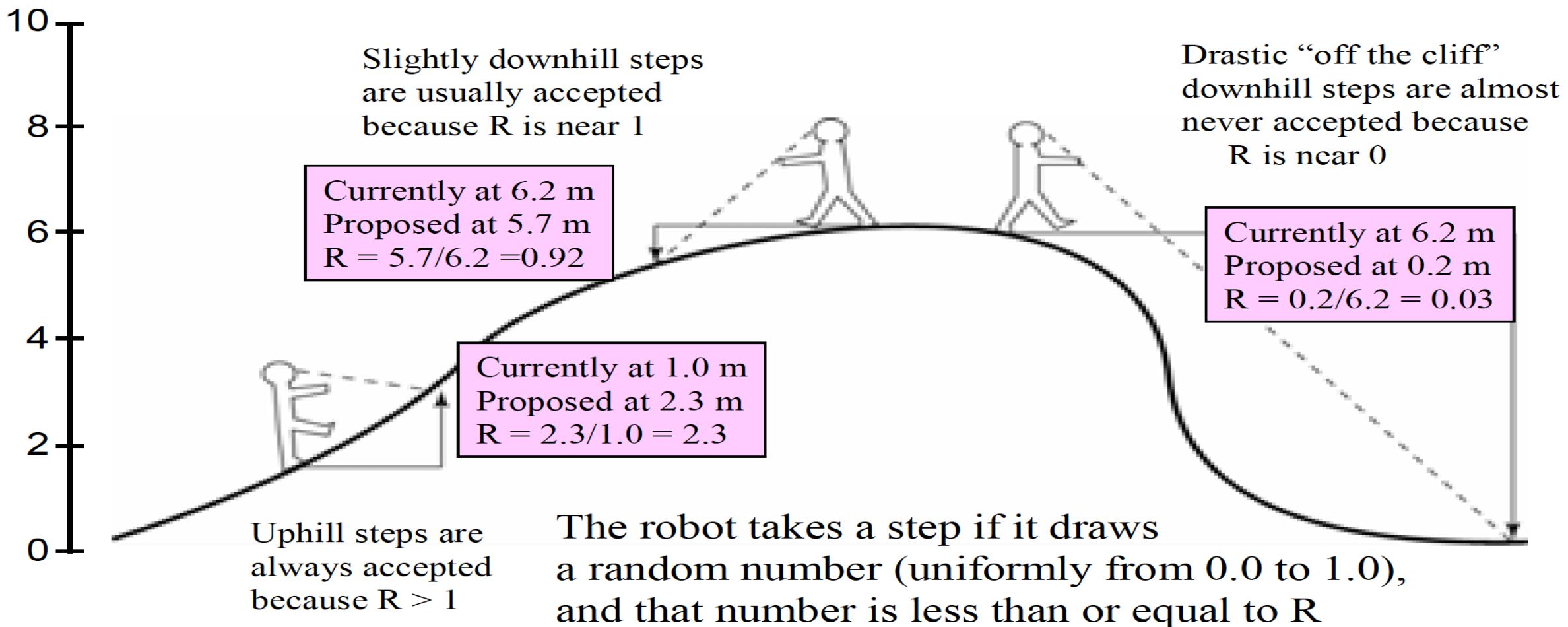
How the MCMC works

- Start somewhere
 - That “somewhere” will have a likelihood associated with it
 - Not the optimized, maximum likelihood
- Randomly propose a new state
 - If the new state has a better likelihood, the chain goes there

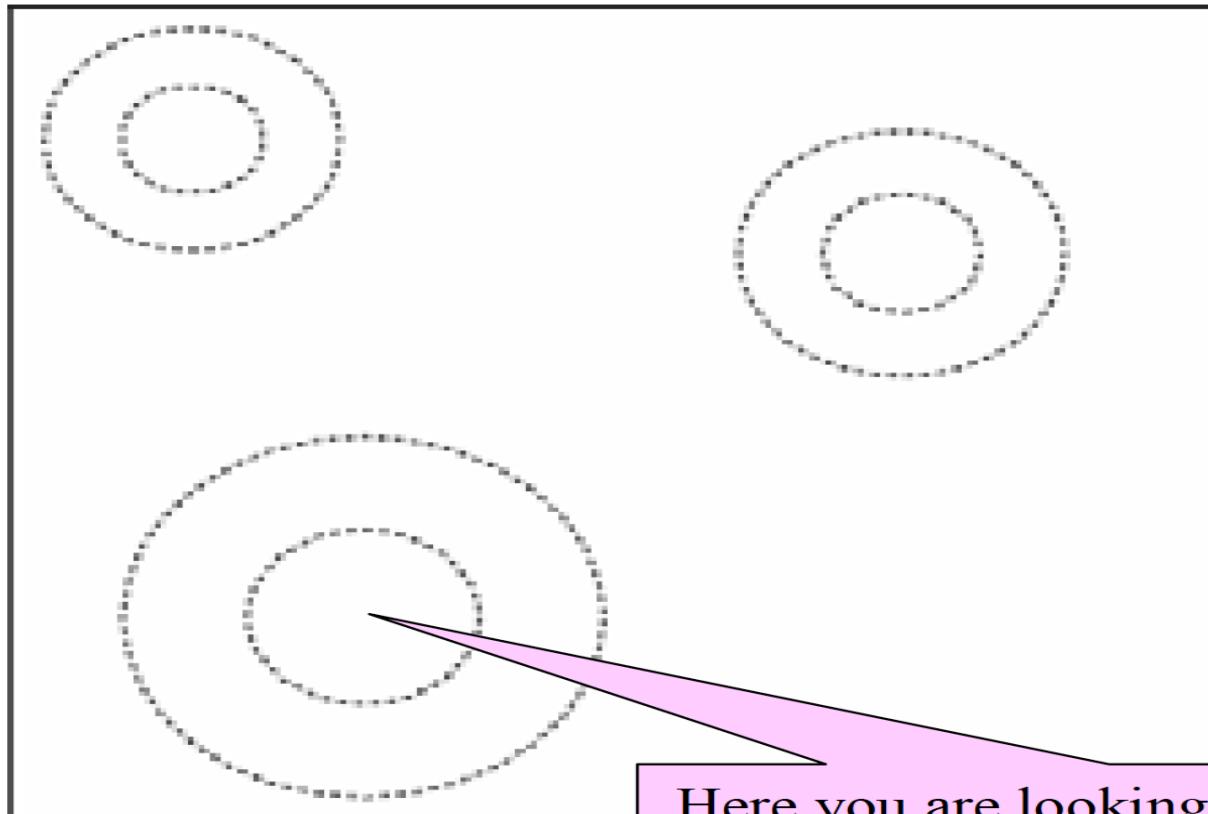
MCMC robot's rules



(Actual) MCMC robot rules



What MCRobot can teach us about Markov chain Monte Carlo



Here you are looking down from above at
one of the three bivariate normal hills

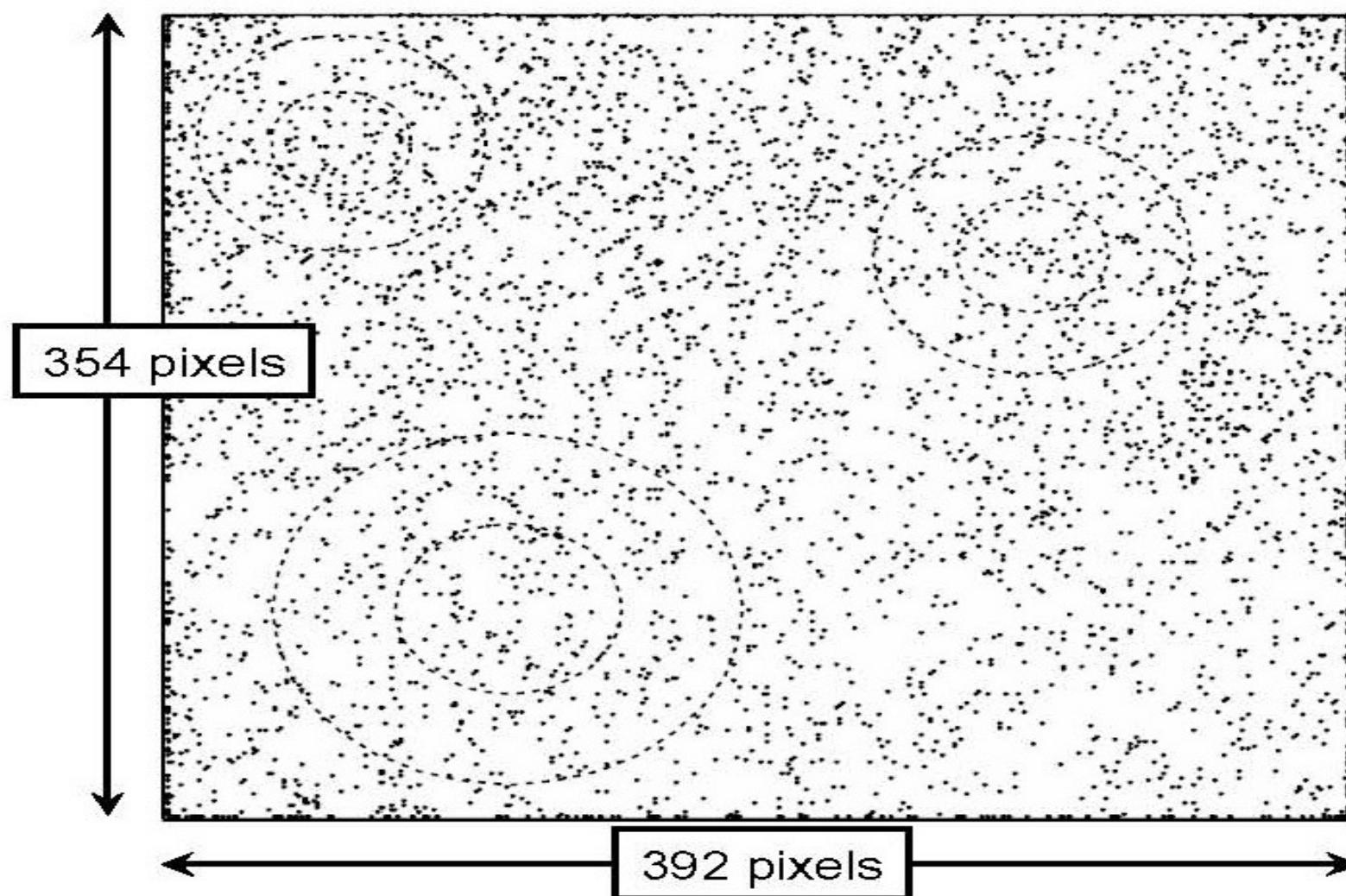
Posterior distribution:

- equal mixture of 3 bivariate normal “hills”
- inner contours: 50%
- outer contours: 95%

Proposal scheme:

- random direction
- gamma-distributed step length
- reflection at edges

Pure random walk



Proposal scheme:

- random direction
- gamma-distributed step length (mean 45 pixels, s.d. 40 pixels)
- reflection at edges

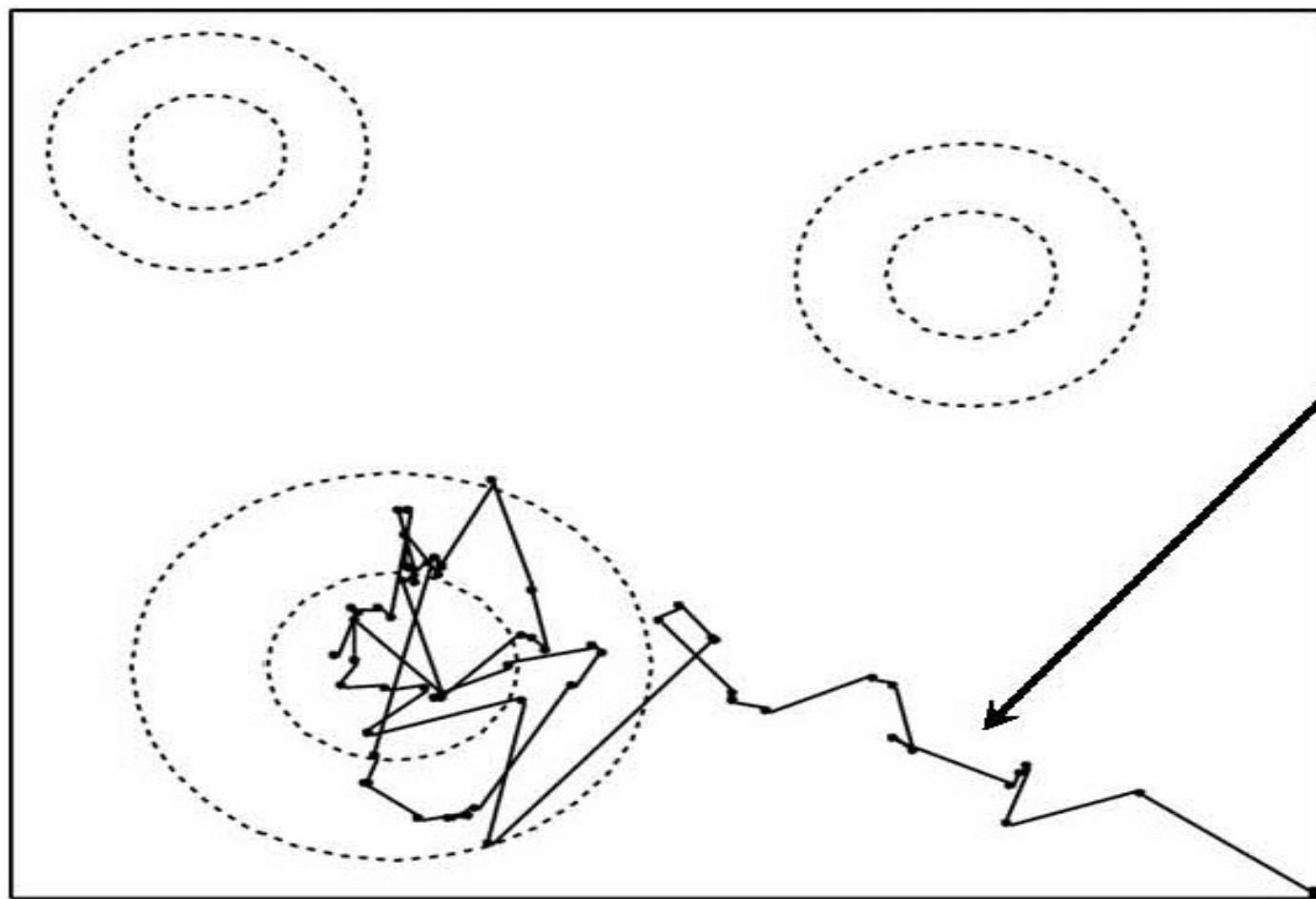
Target distribution:

- equal mixture of 3 bivariate normal "hills"
- inner contours: 50%
- outer contours: 95%

In this case, the robot is accepting every step

5000 steps shown

Burn-in



Robot is now following the rules and thus quickly finds one of the three hills.

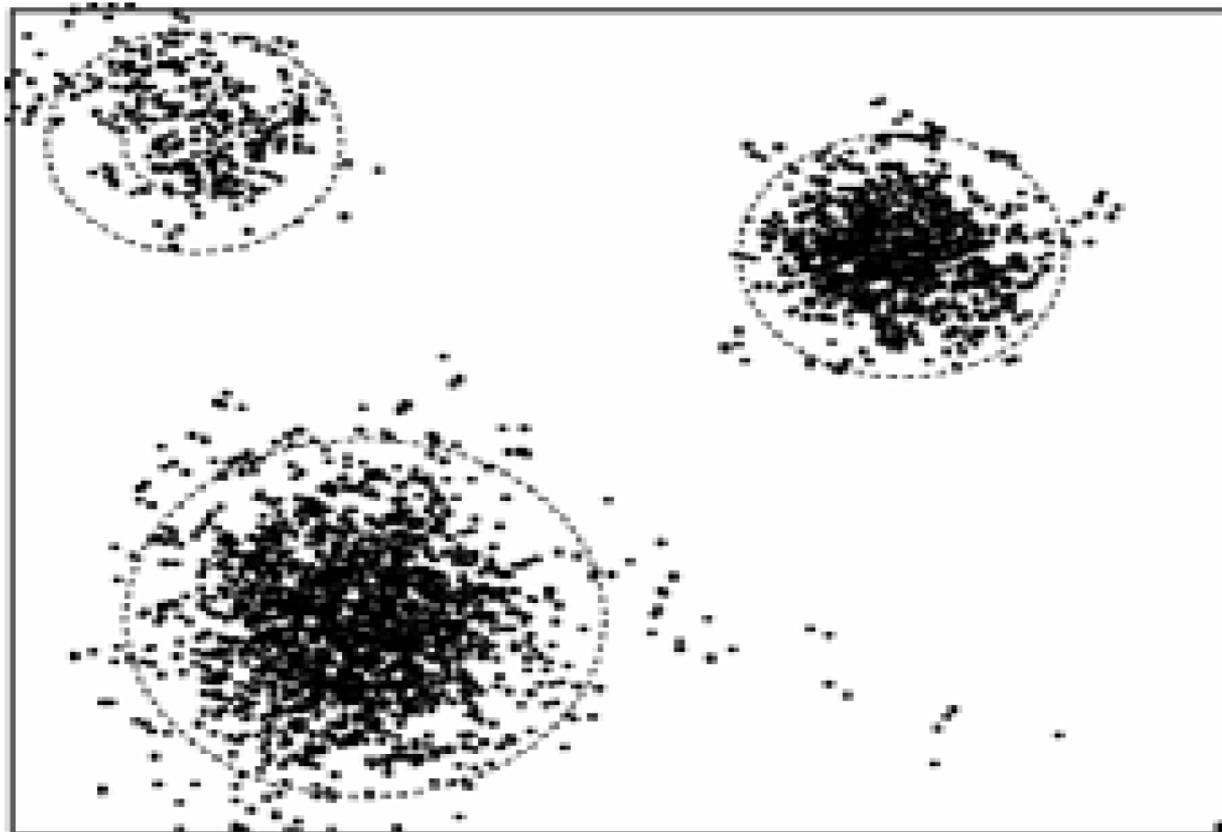
Note that first few steps are not at all representative of the distribution.

100 steps taken

Starting point

Target distribution approximation

5000 steps taken



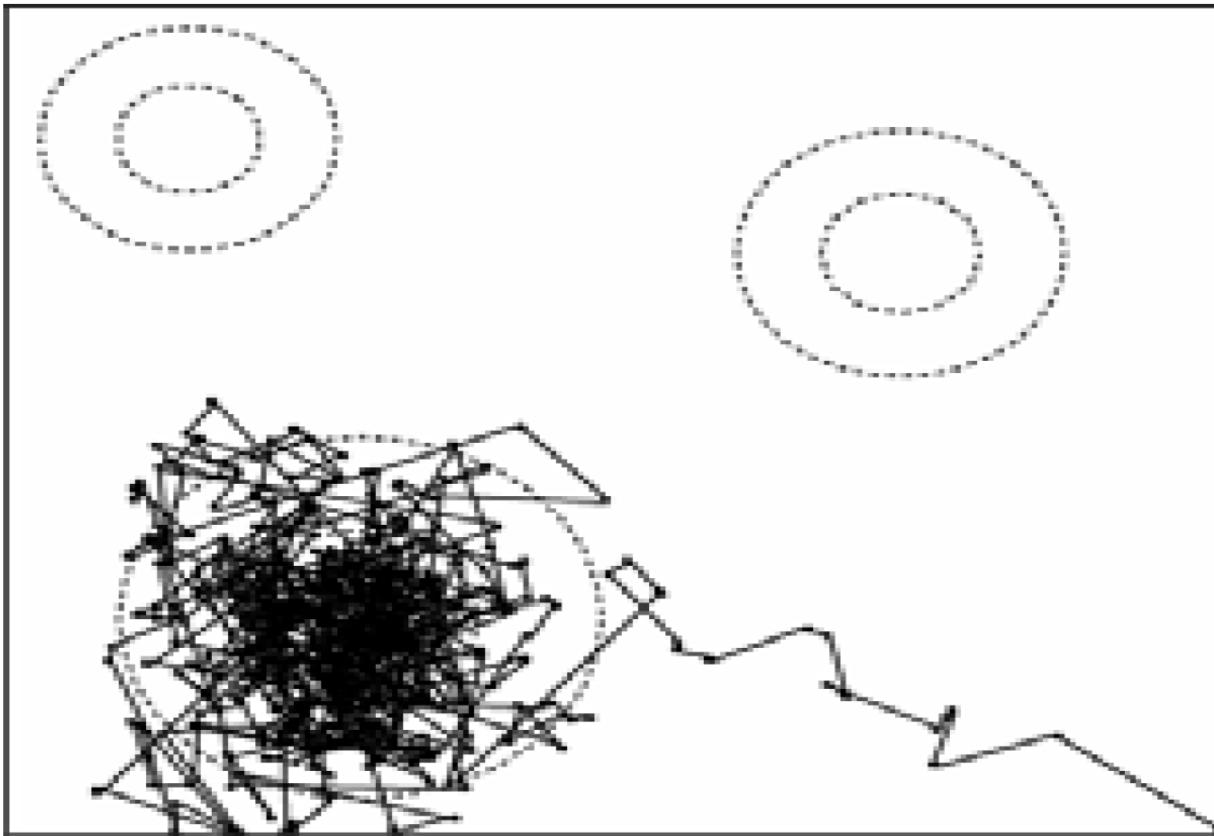
How good is the MCMC approximation?

- 51.2% of points are inside inner contours (cf. 50% actual)
- 93.6% of points are inside outer contours (cf. 95% actual)

Approximation gets better the longer the chain is allowed to run.

Just how long is a long run?

1000 steps taken



What would you conclude about the target distribution had you stopped the robot at this point?

The way to avoid this mistake is to perform **several runs**, each one beginning from a different randomly-chosen starting point.

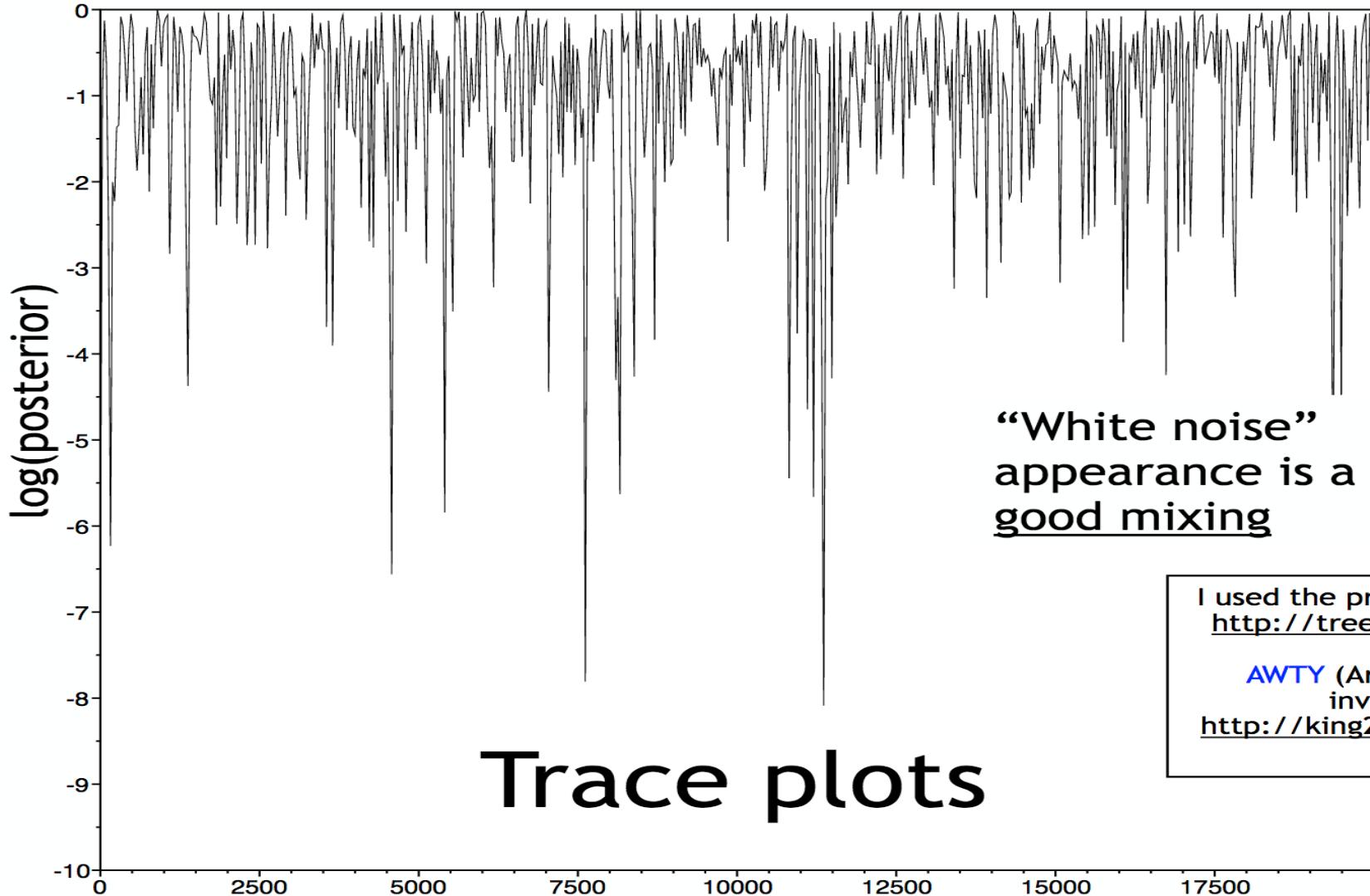
Results different among runs? Probably none of them were long enough!

Bayesian phylogenetics

- We have rules for how we change parameters and trees
- During each step, we change one of the parameters according to its rule
 - Can be a model parameter, e.g. alpha in the JC model
 - Can be a branch length
 - Can be the topology
 - E.g. using tree rearrangements that we talked about

MCMC, in short:

- Start with **random tree** and arbitrary **initial values for branch lengths and model parameters**
- Each generation consists of one of these (chosen at random):
 - **Propose a new tree** (e.g. NNI) and either accept or reject the move
 - **Propose** (and either accept or reject) a new model **parameter value**
- Every k generations save tree topology, branch lengths and all model parameters (i.e. sample the chain)
- After n generations, summarize sample using histograms, means, credible intervals, etc.

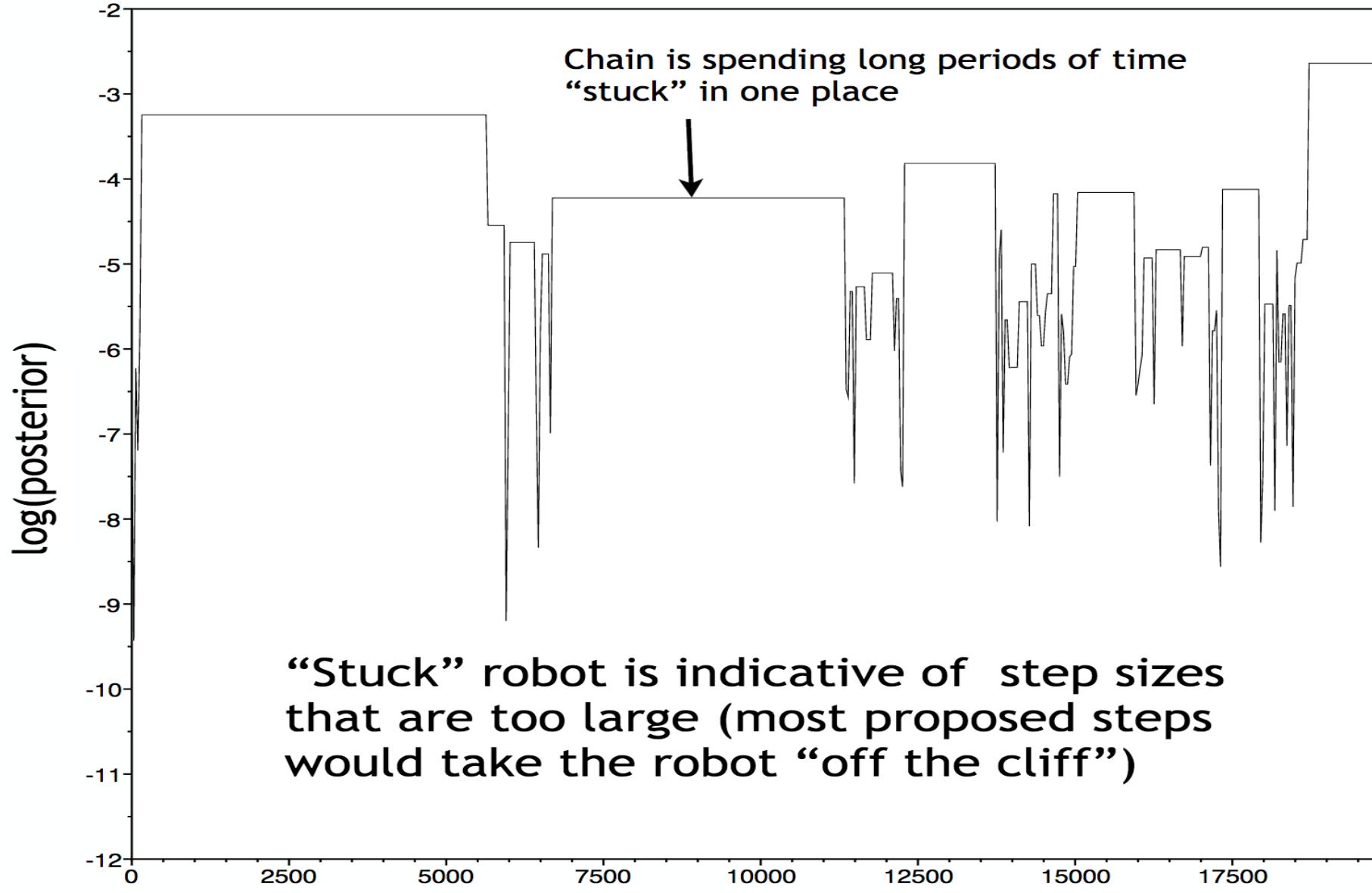


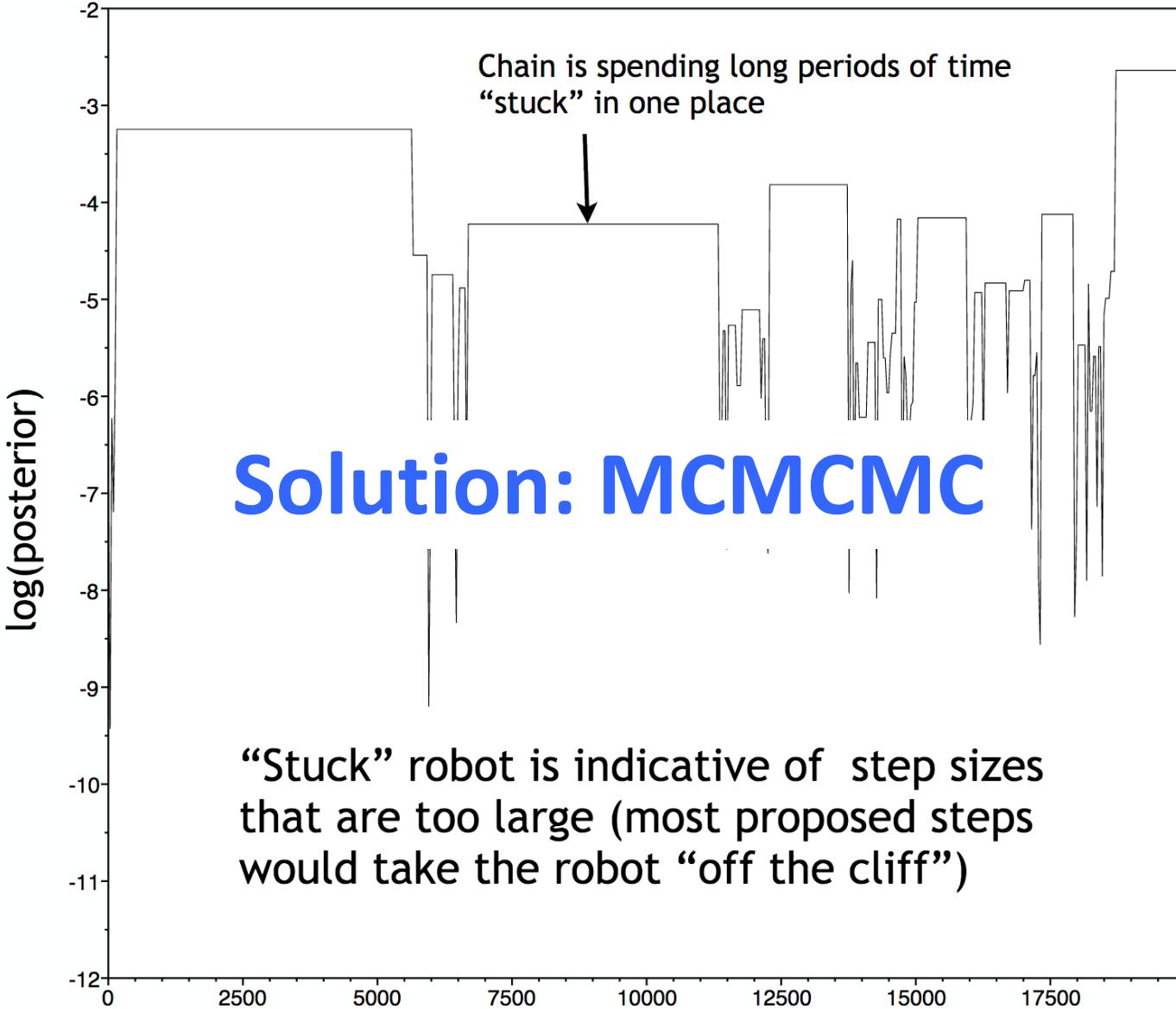
I used the program [Tracer](#) to create this plot:
<http://tree.bio.ed.ac.uk/software/tracer/>

[AWTY](#) (Are We There Yet?) is useful for
investigating convergence:
[http://king2.scs.fsu.edu/CEBProjects/awty/
awty_start.php](http://king2.scs.fsu.edu/CEBProjects/awty/awty_start.php)

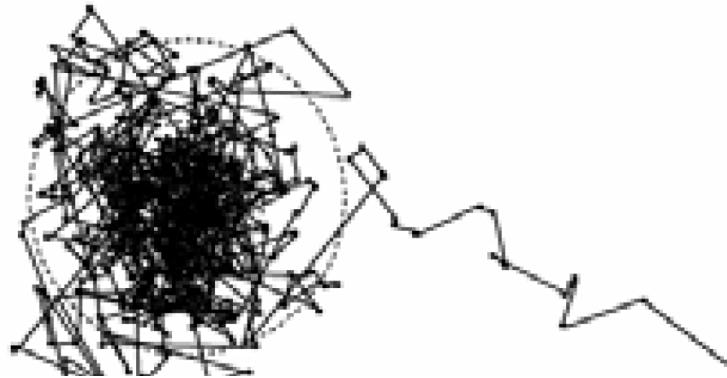
Trace plots







Metropolis-coupled Markov chain Monte Carlo (MCMCMC, or MC³)



- MC³ involves running **several chains simultaneously** (one “cold” and several “heated”)
- The cold chain is the one that counts, the heated chains are “scouts”
- Chain is heated by raising densities to a power less than 1.0 (values closer to 0.0 are warmer)

Cold vs. heated landscapes

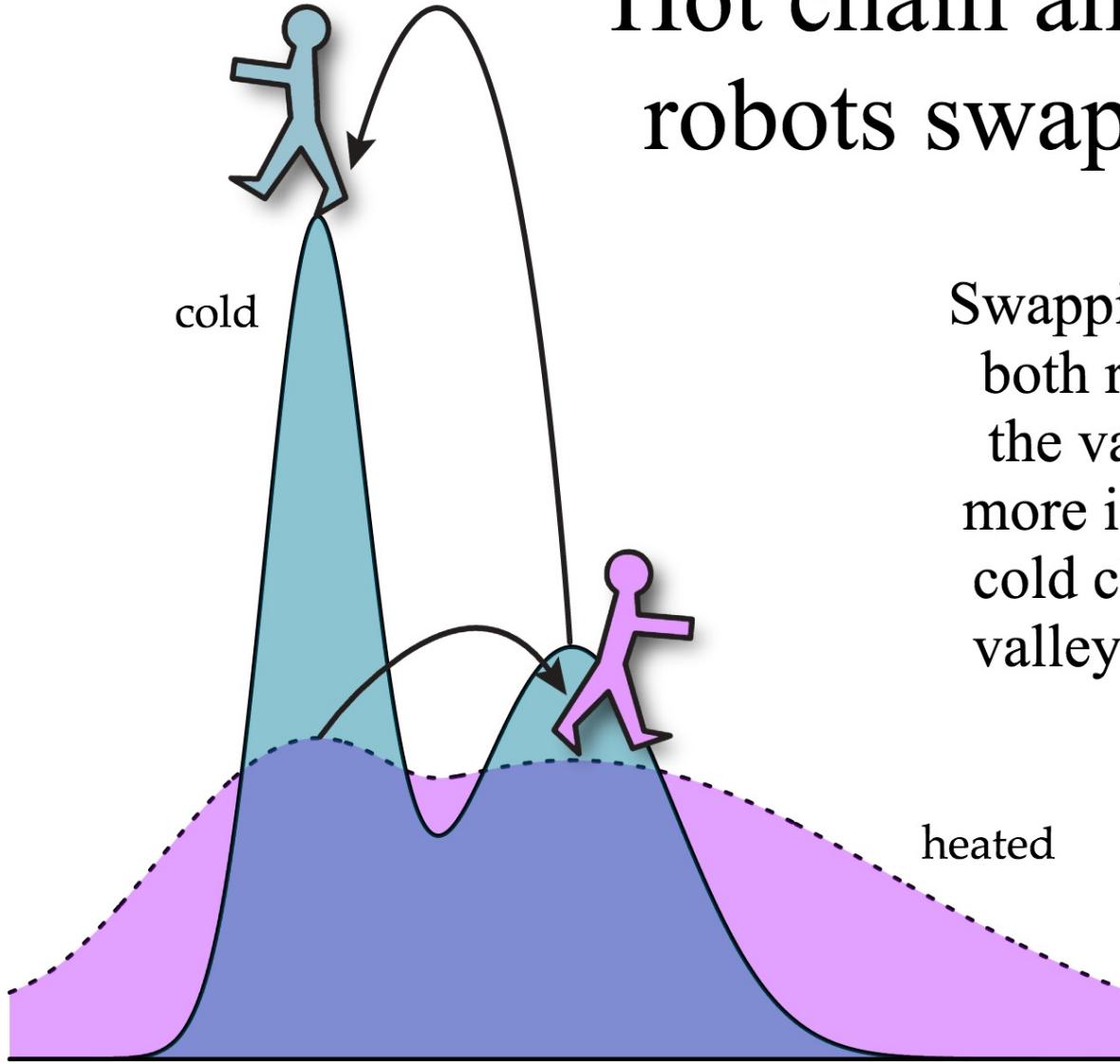


Cold landscape: note peaks separated by deep valleys



Heated landscape: note shallow (easy to cross) valleys

Hot chain and cold chain robots swapping places



Swapping places means both robots can cross the valley, but this is more important for the cold chain because its valley is much deeper

How does MCMC work?

A simple summary ☺

1/2

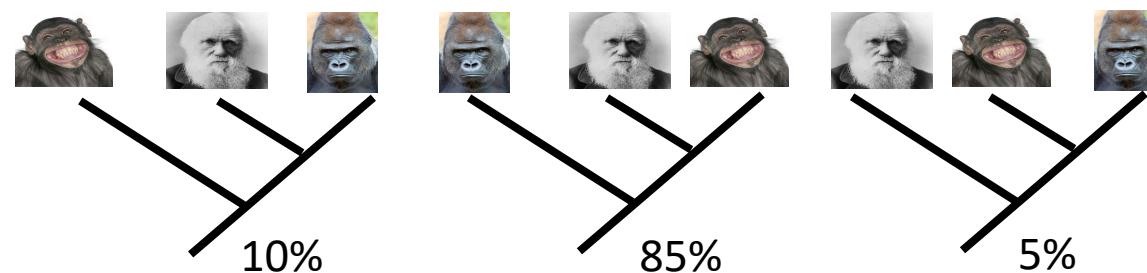
- A Markov chain generates a series of random variables
- The probability distribution of future states depends only on the current state (Markov property)
- Phylogenetic inference starts with a randomly generated tree with branch lengths
- Next step is to generate a new tree, based on the previous tree e.g. using tree rearrangements (NNI, SPR, TBR) or changing branch lengths -> this is a **proposal**
- The **proposal is accepted or rejected** given a probability based on the Metropolis-Hastings algorithm
 - In practice, it's accepted if it has a better likelihood
- If it is accepted, it becomes the new current state and a new proposal is made

How does MCMC work?

2/2

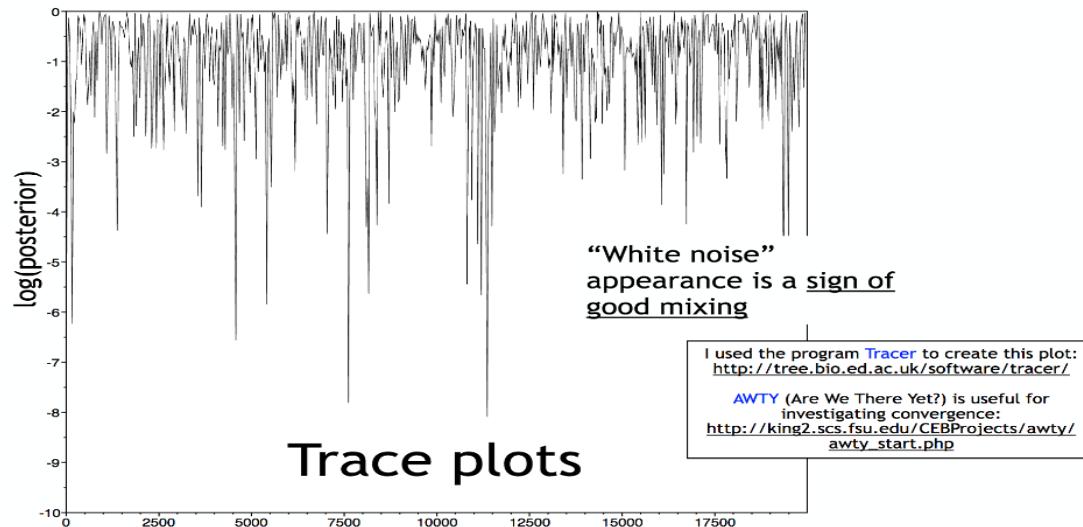
A simple summary ☺

- Running a Markov chain relatively quickly finds better trees
- After a while no better trees can be found and all sampled trees are close to the optimum – “stationary distribution”
- The number of times the tree is visited by the chain – interpreted as posterior probability of that tree

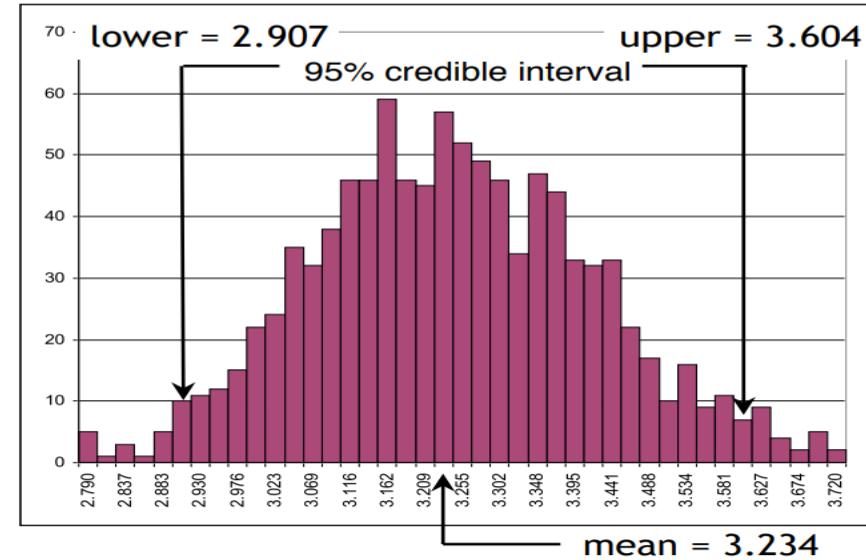


Looking at the results

- Graphically



Marginal Posterior Distribution of κ



Paul O. Lewis (2014 Woods Hole Molecular Evolution Workshop)

Data from Lewis, L., and Flechtner, V. 2002. Taxon 51: 443-451.

70

- **Statistically, by computing Effective Sample Size (ESS, minimum should be 200 or more)**
- **Checking convergence of each run, and all runs together**

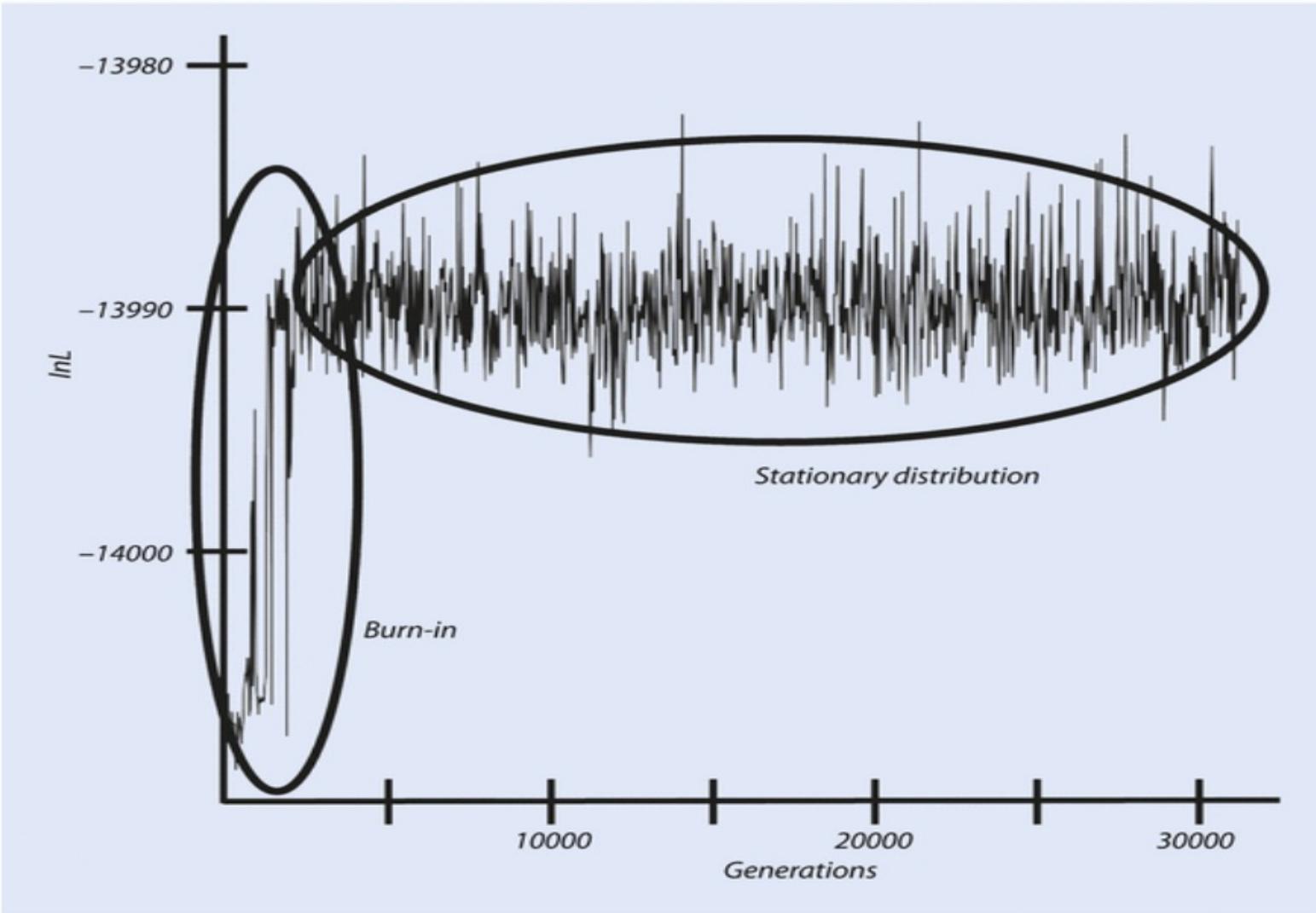
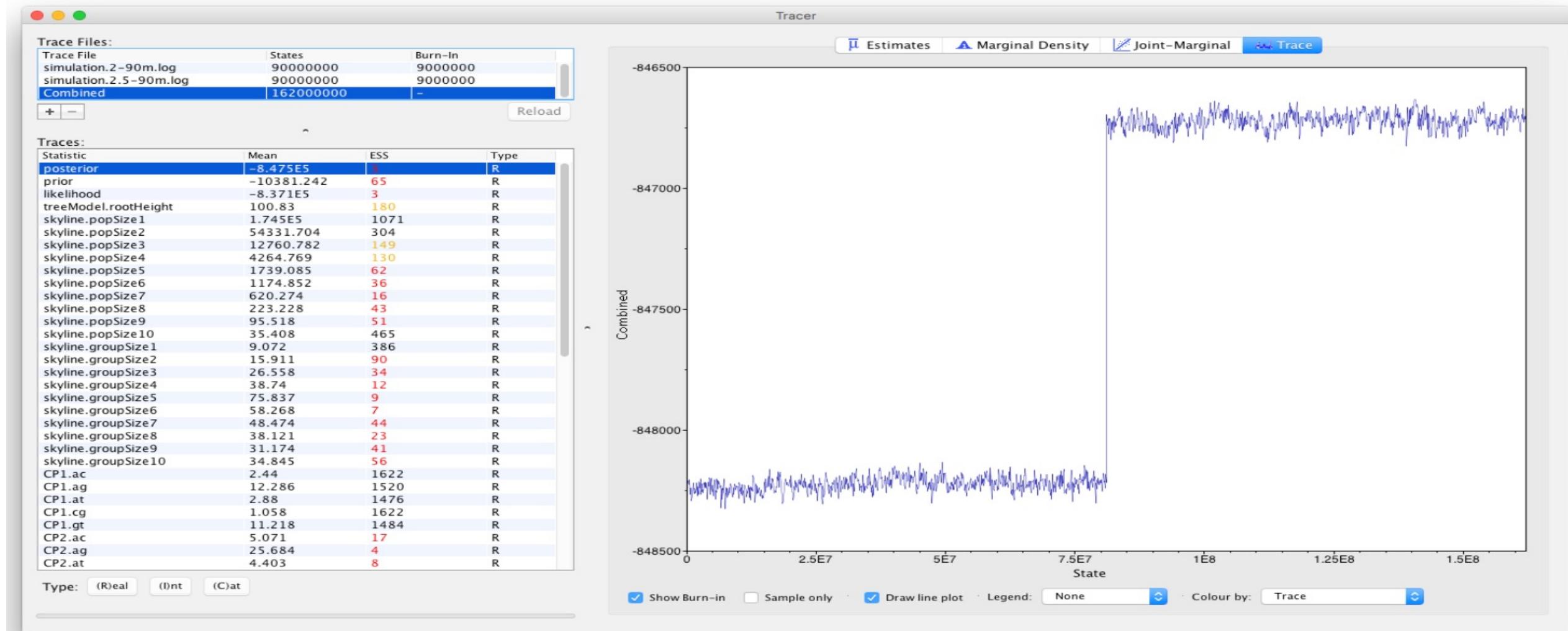


Fig. 8.12 Likelihood scores of a MCMC run plotted against generations. Once stationarity is achieved, trees from this distribution are sampled by discarding all other trees as burn-in. A majority-rule consensus of sampled trees will provide posterior probabilities for every node

Lack of convergence



Summarizing the results

- ▶ **Autocorrelation:** between values that are sampled one after another, so need to sample values at a lower frequency – e.g. every 1000 steps
- ▶ MCMC easily runs over millions of states

=> **Synthesize/summarize the parameters we are interested in**

By computing the marginal posterior distribution of these parameters

- mean, median or variance
- 95% credibility interval

By identifying one or more “best” topologies

e.g. the splits most frequently identified

The number of times a clade in the tree is accepted during the MCMC defines the posterior probability of the clade, and therefore indicates the support for the node

In summary

- Statistical frameworks allow us to estimate parameters for our models, including tree topology
- Maximum likelihood gives us point estimates
 - estimates that maximize the likelihood of the data given the models
- Bayesian inference allows us to explore the error around our estimates
 - posterior probabilities of the models given the data