# BAYESIAN MCMC

# aims

To come away with an intuitive* understanding of:

(1) maximum likelihood
(2) Bayes Theorem
(3) what they mean for phylogenetic inference

(4) a basis for understanding how** they work

*as little mathematics as possible
**requires equations

# We find the set of trees that maximise the Lh

$$P(X|\theta) = Lh$$

X = data
θ = model parameters

{tree topology, branch lengths, node values}
{mathematical description of change}

Take a deep breath

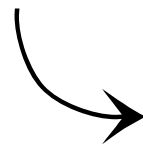# Calculating the likelihood
# X = observed data

A: Matrix

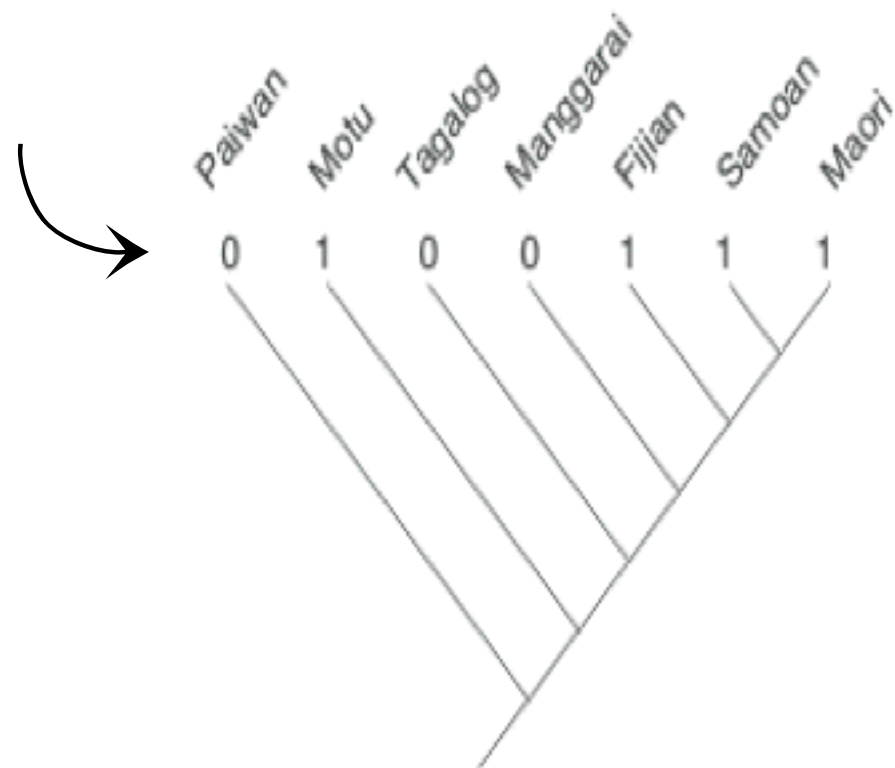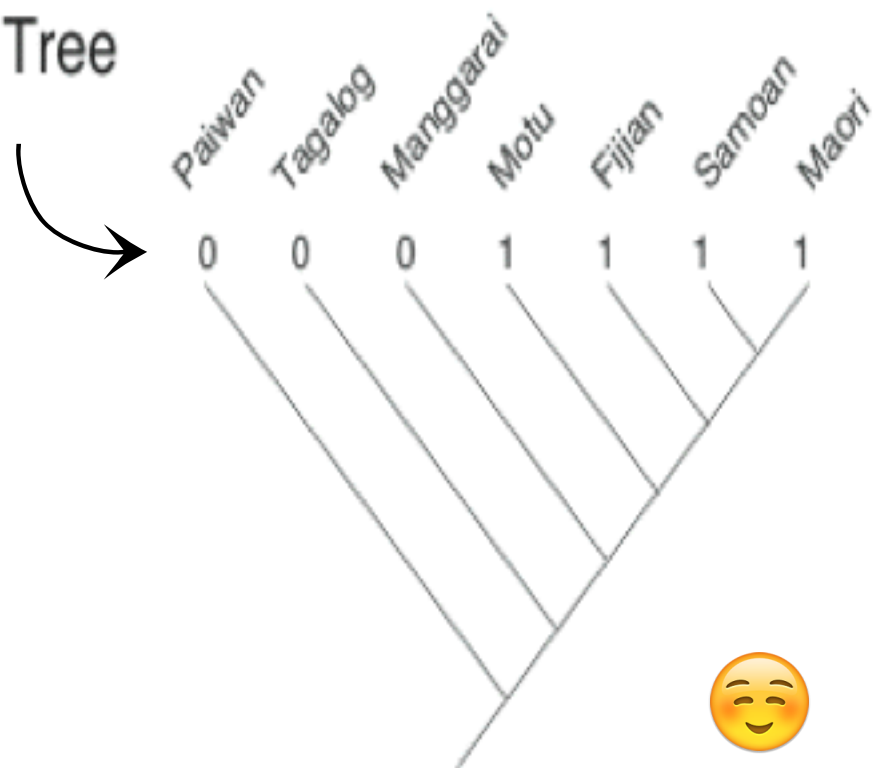| | a | b | c | d | e | f | g | h | I | j | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Paiwan | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | ... |
| Tagalog | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | ... |
| Manggarai | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | ... |
| Motu | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | ... |
| Fijian | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | ... |
| Samoan | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | ... |
| Maori | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | ... |

# Calculating the site likelihood
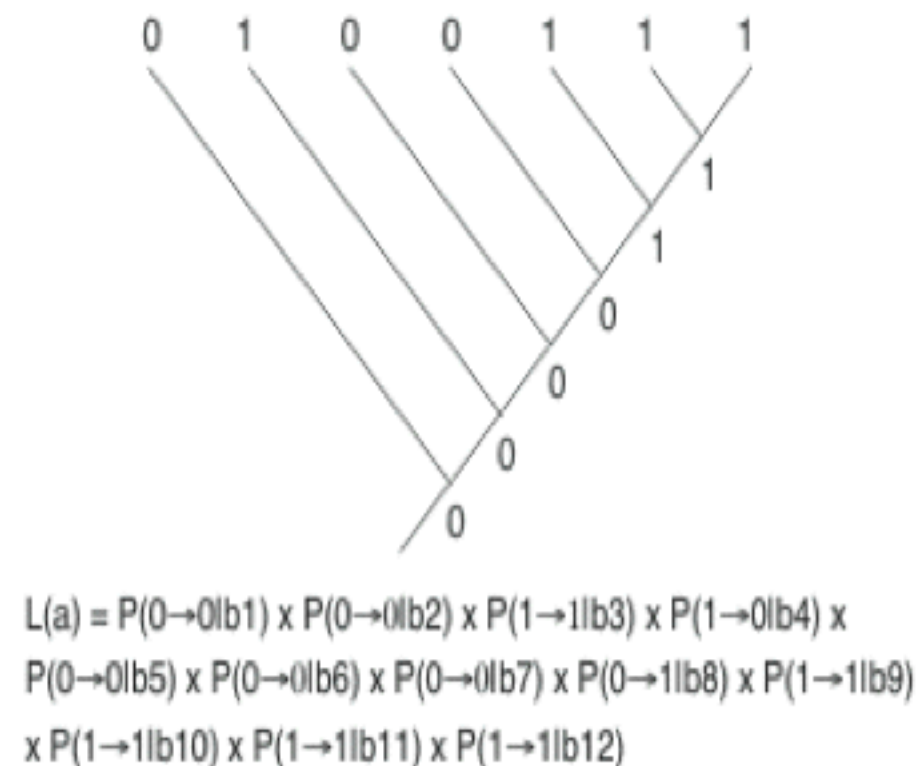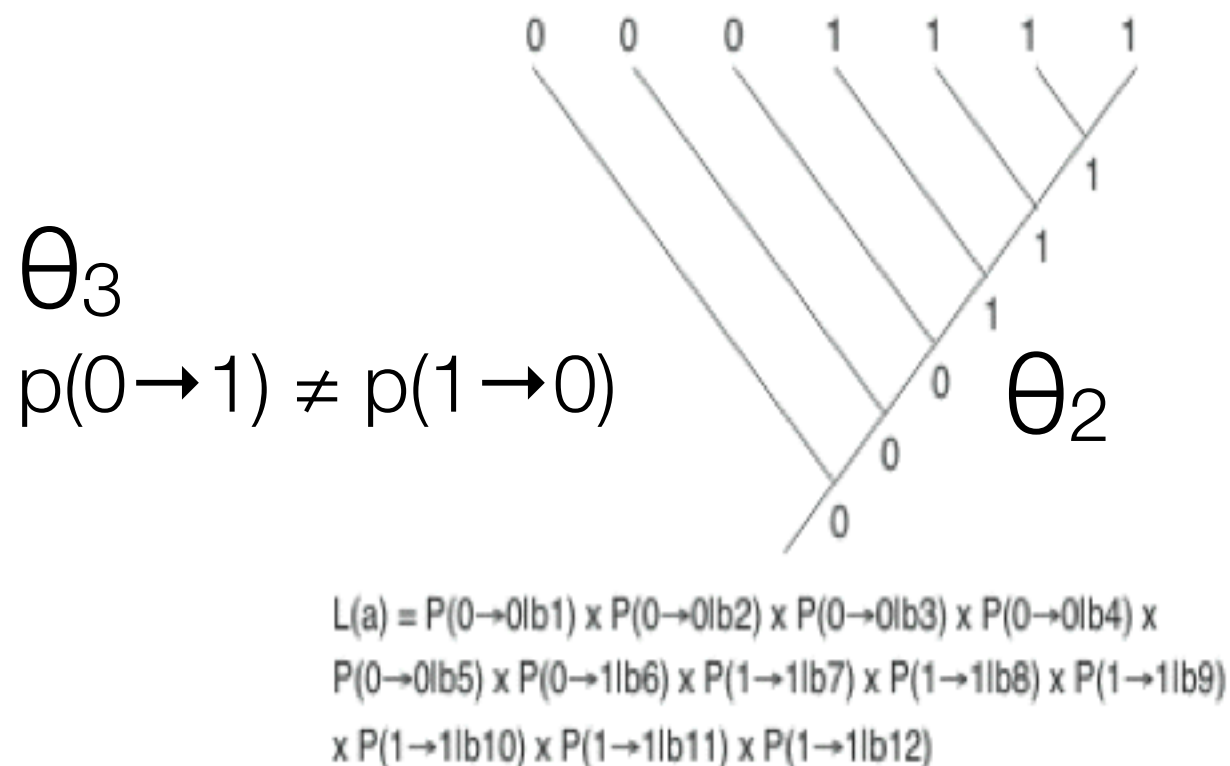# $\theta_1$ = trees (some hypotheses about history)

# Calculating the site likelihood
# $\theta_2$ = node states and $\theta_3$ = model of change

We combine a model of character change with some node values



C: Ancestral States

$\theta_3$
$p(0 \to 1) \neq p(1 \to 0)$

$\theta_2$

L(a) = P(0→0|b1) x P(0→0|b2) x P(0→0|b3) x P(0→0|b4) x
P(0→0|b5) x P(0→1|b6) x P(1→1|b7) x P(1→1|b8) x P(1→1|b9)
x P(1→1|b10) x P(1→1|b11) x P(1→1|b12)

L(a) = P(0→0|b1) x P(0→0|b2) x P(1→1|b3) x P(1→0|b4) x
P(0→0|b5) x P(0→0|b6) x P(0→0|b7) x P(0→1|b8) x P(1→1|b9)
x P(1→1|b10) x P(1→1|b11) x P(1→1|b12)

IN WORDS: the likelihood of character (a) is equal to the probability of (a) staying in state 0 along branch 1, multiplied by the probability of (a) staying in state 0 along branch 2 (...), multiplied by the probability of (a) changing to state 1 on branch 6, multiplied by the probability of it staying in state 0 on branch 7 ...

from Greenhill & Gray 2009

# Calculating the site likelihood
# Lh = node values over all possible values ($\Pi\theta_2$)



D: Site Likelihood

Site Likelihood(a) = P(reconstruction 1) ...x... P(reconstruction 5) ... x ... P(reconstruction n)

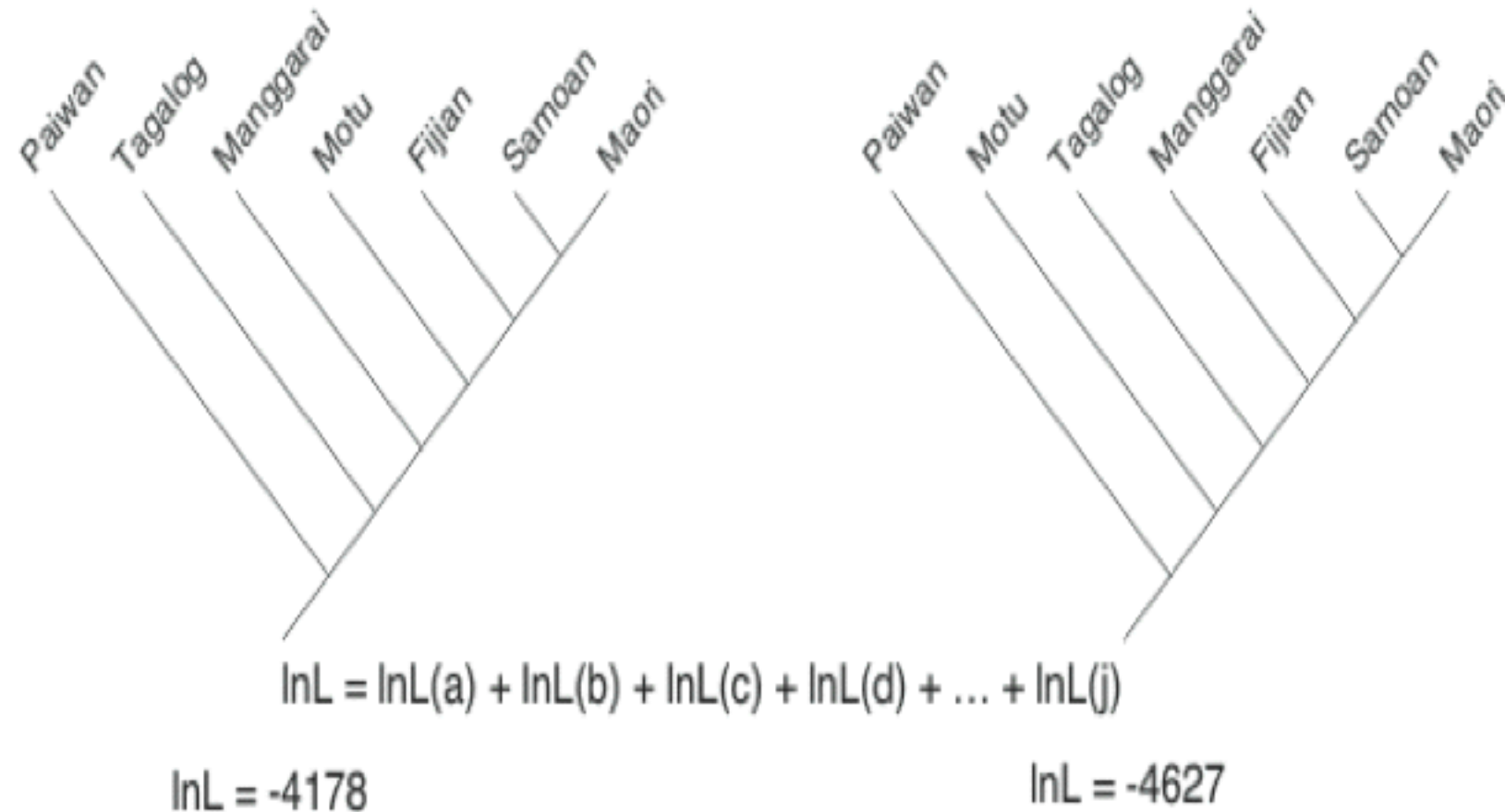**...** multiplied over all the possible node values for the character.

# Calculating the tree likelihood
# Lh = Σ site likelihoods

E: Tree Likelihood

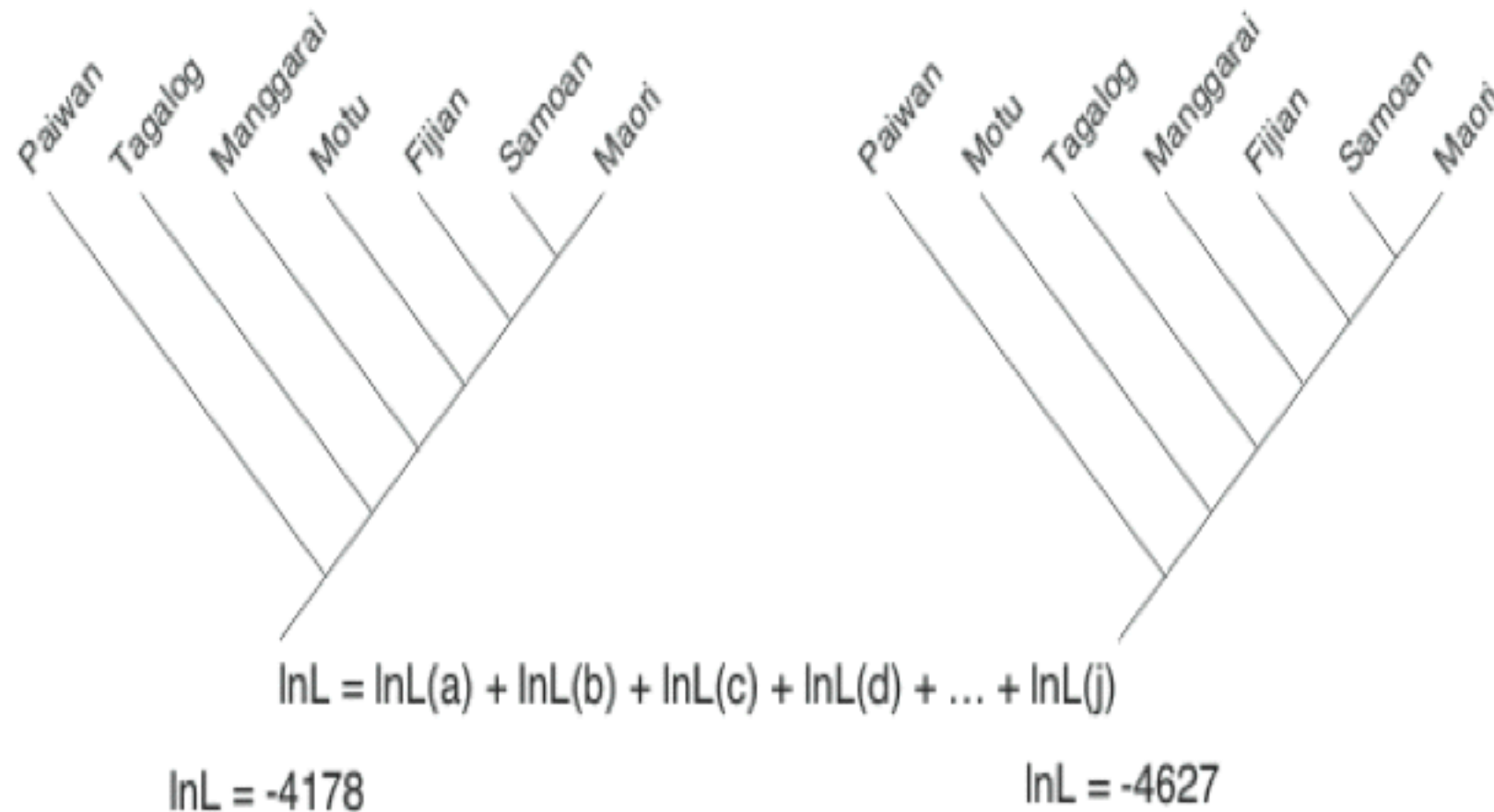Paiwan Tagalog Manggarai Motu Fijian Samoan Maori

Paiwan Motu Tagalog Manggarai Fijian Samoan Maori

$$lnL = lnL(a) + lnL(b) + lnL(c) + lnL(d) + ... + lnL(j)$$

lnL = -4178

lnL = -4627

IN WORDS: The likelihood of the tree is the sum of the likelihood of each site on that tree.

# Calculating the maximum likelihood
# Find the tree with the best Lh

E: Tree Likelihood



$$\ln L = \ln L(a) + \ln L(b) + \ln L(c) + \ln L(d) + \ldots + \ln L(j)$$

InL = -4178

InL = -4627

The higher the likelihood (closer to zero) the more we prefer a tree. We retain that tree and tweak it in an iterative fashion.

We infer character state changes, so innovations and retentions are incorporated in the model.

from Greenhill & Gray 2009

# Finding the maximum likelihood is computationally expensive

# the good and the not-so-good of maximum likelihood

- desirable statistical properties
- explicit expression of model
- model-testing framework (LRT)
- lots of data > converge on MLE

- non-intuitive
- computationally expensive
- how do we integrate over trees?
- how to account for uncertainty

# the Bayesian approach

- explicit expression of model
- model-testing framework
- retains advantages of ML
- computationally efficient
- can integrate over trees
- can account for uncertainty
- best for linguistic/cultural data

- intuitive statistical reasoning?

Huelsenbeck & Ronquist 2009 / Ronquist 2004

reasoning about probability and updating our reasoning as new information becomes available

# talking in maths*

P
probability

t
time

Σ
sum
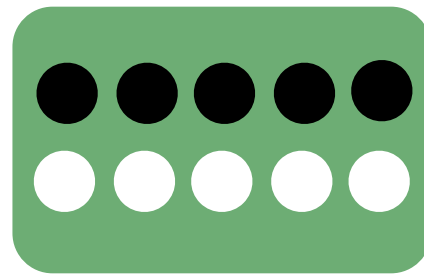
θ
"model"

q
"rate"

X
data

∫
integral

|
conditional

# probability

"forward" probability

M = 10 black, 5 white

X = 5 black then 5 white

$p(\mathbf{X}|M) = 10/15 * 9/15 * ... * 1/6 = 0.00033$
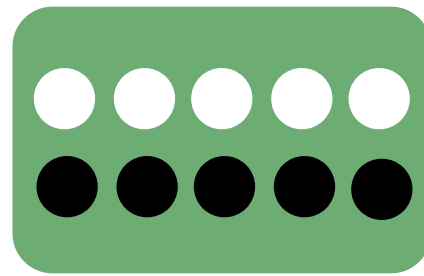
i.e. what is the probability of these balls (data) given this model (urn)?
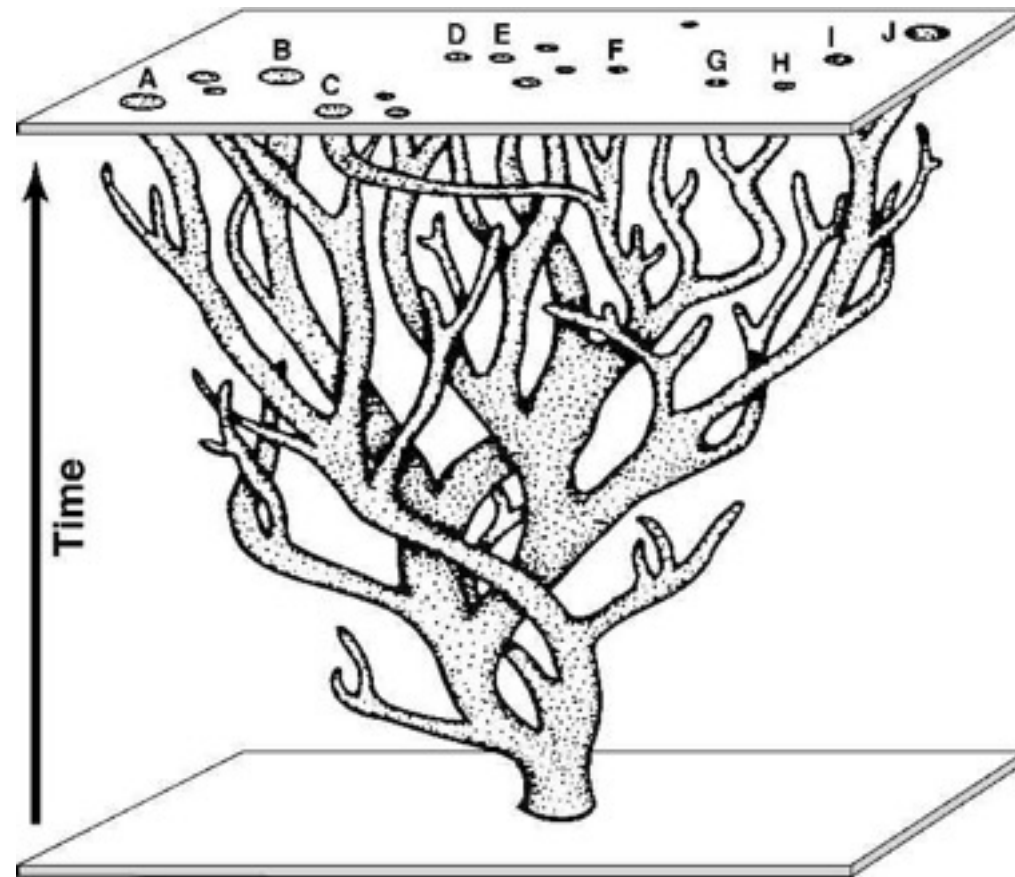
**"converse" probability**

M = ?

X = 5 black, 5 white

How do we find p(M|X) ???
Impossible unless we know something about M!

i.e. given these balls (data) what can we infer about what's in the urn (model)?

"converse" probability



The historical sciences regularly deal with this situation.

i.e. given these balls (data) what can we infer about what's in the urn (model)?

$\mathfrak{X}$ are the Data

$\Theta$ the model Parameters

prior       likelihood

$$f(\theta \mid X) = \frac{f(\theta)\,f(X \mid \theta)}{\int f(\theta)\,f(X \mid \theta)\,d\theta}$$

posterior

normalizing constant

T. Bayes.

# Take another deep breath

# Bayes Theorem 2

conditional probability  $p(X|M) = \dfrac{p(XM)}{p(M)}$

joint probability  $p(M) \times p(X|M) = p(XM)$

similarly  $p(X) \times p(M|X) = p(XM)$

remember the balls (X) and urn (M) situation:  $p\left(\begin{array}{c|c} \blacksquare & \blacksquare \end{array}\right)$

# Bayes Theorem 3

## joint probability

$$p(M) \times p(X|M) = p(XM)$$

similarly

$$p(X) \times p(M|X) = p(XM)$$

$$p(M|X) = \frac{p(XM)}{p(X)}$$

$$p(M|X) = \frac{p(M) \times p(X|M)}{p(X)}$$

## Bayes Theorem

$$p(\theta|X) = \frac{p(X|\theta)\, p(\theta)}{p(X)}$$

Bayes Theorem

$$p(\theta|X) = \frac{p(X|\theta)\,p(\theta)}{p(X)}$$

$$p(\theta|X) = \frac{p(X|\theta)\,p(\theta)}{\sum \{(p(X|\theta_1)p(\theta_1)) + (p(X|\theta_2)p(\theta_2)) + ...\}}$$

$$p(\theta|X) = \frac{p(X|\theta)\,p(\theta)}{\int p(X) \text{ over all } \theta}$$

$$p(\text{you}) = \text{argh!} = 1.0$$

Bayes Theorem

$$p(\theta|X) = \frac{p(X|\theta)\, p(\theta)}{p(X)}$$

the probability of the data given the hypothesis

the probability of the hypothesis

$$\text{posterior probability} = \frac{\text{likelihood of data} \times \text{prior}}{\text{marginal likelihood}}$$

the probability of the hypothesis, given the data

the unconditional (over all hypotheses) probability of the data

calculating the Lh(site) is **analytically intractable**

for 100 taxa, there are $4.02 \times 10^{59}$
possible ancestral state configurations for one site

$4.02 \times 100,000,000,000,000,000,000,000,000,000,000,000,000,000,000,000,000,000,000,000,000!$

we need to **sample** this parameter space

random sampling is wasteful (huge space, lots of low probability)

MCMC to the rescue!

Geyer 1992

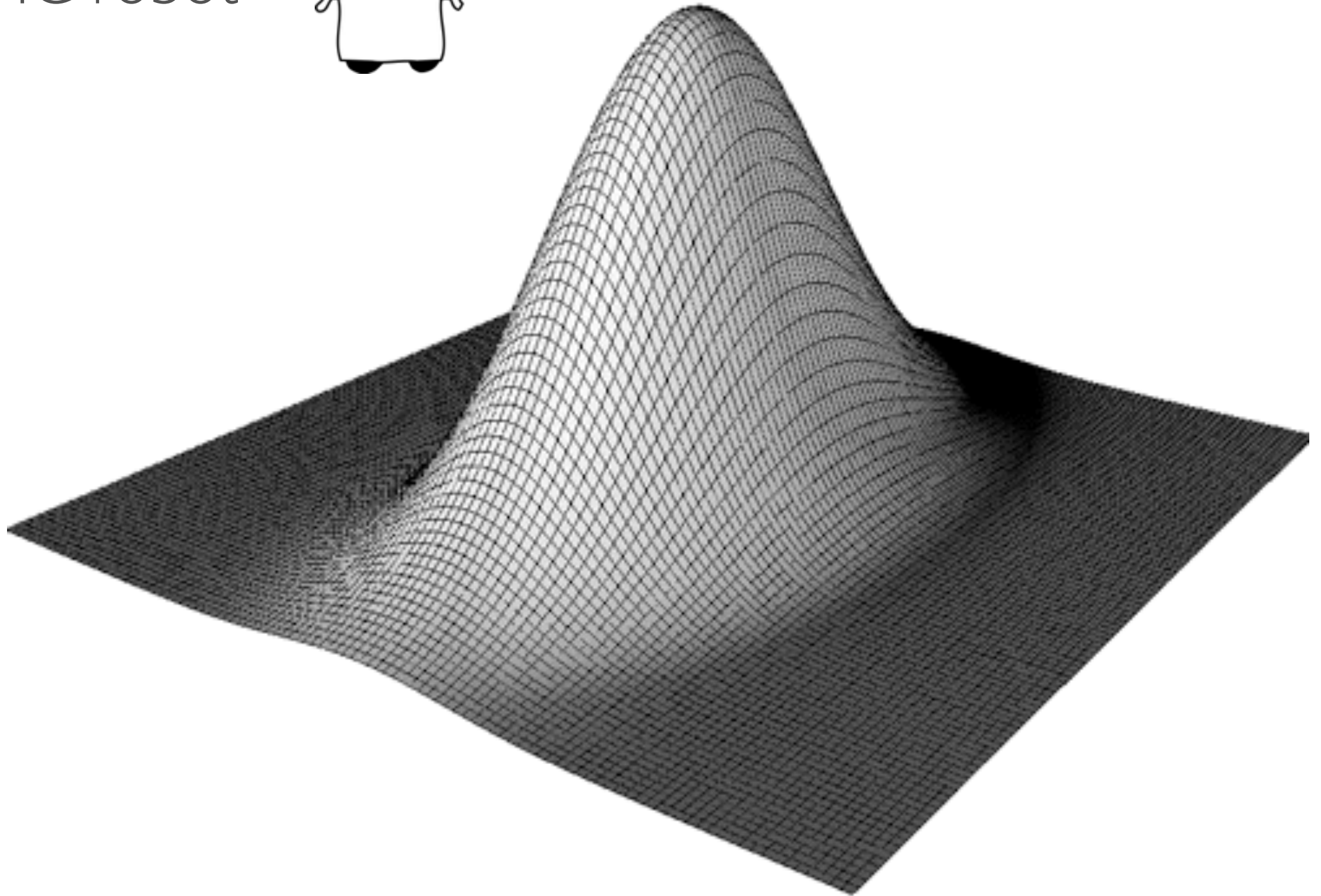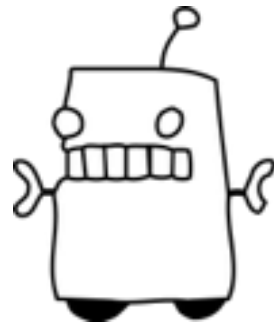the posterior probability distribution as a landscape

"parameter space"

{ topology
branch lengths
node values
transition rates }

{ topology
branch lengths
node values
transition rates }

MCMC robot

MCMC robot

parameter space

transition rates

Markov process

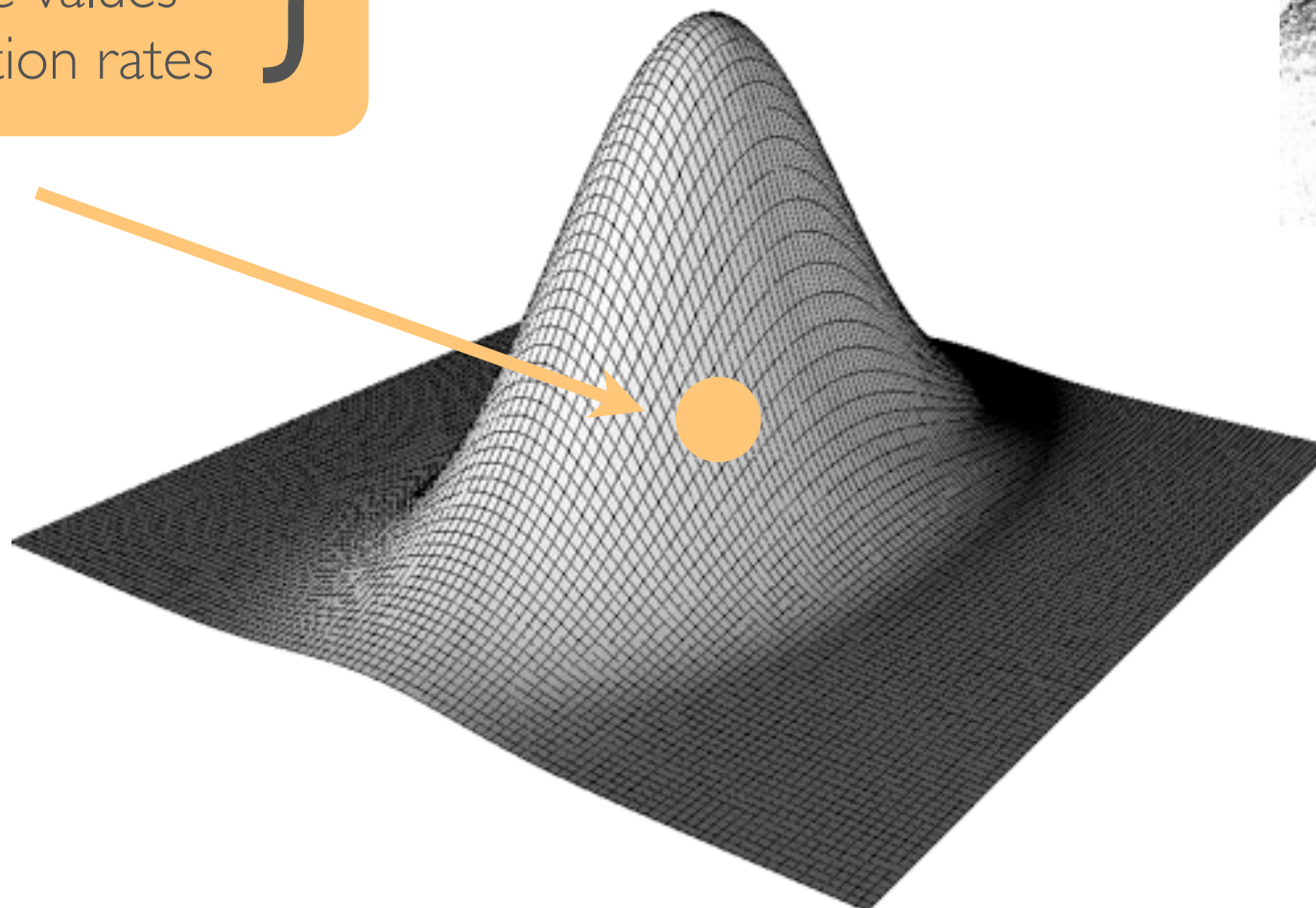posterior probability distribution

Monte Carlo:
repeated random sampling

𝒳 are the Data
Θ the model Parameters
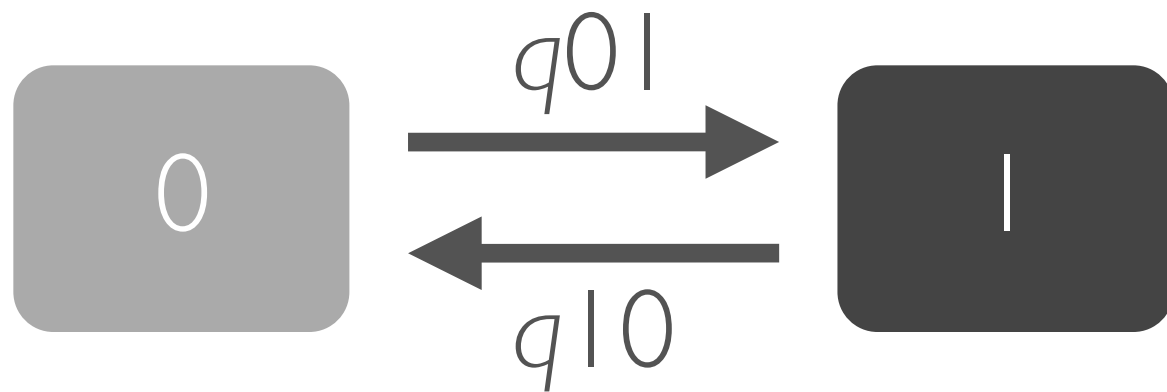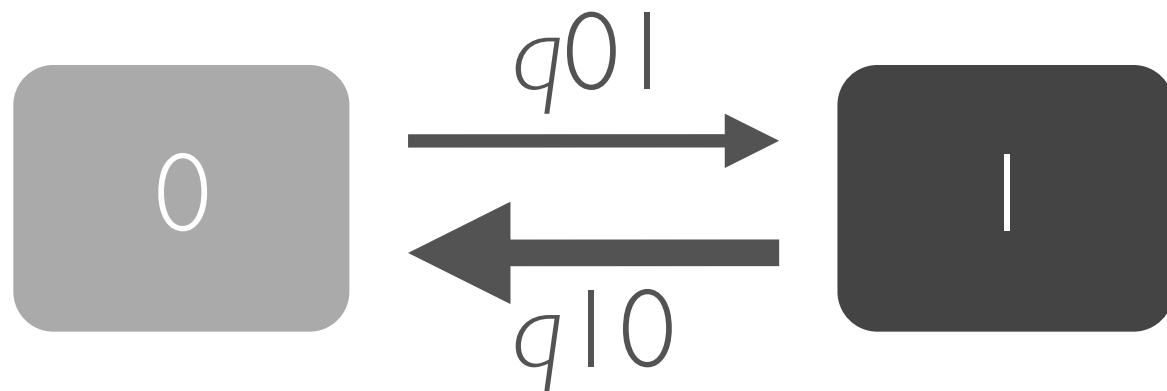
{ topology
branch lengths
node values
transition rates }

J. Bayes.

modelling change in discrete characters 1

$q01 = q10$
a one-parameter model

$q01 \neq q10$
a two-parameter model

$q01 = \alpha$
$q10 = \beta$

we want to estimate these transition rate parameters!

Pagel 1994

# Markov process

"A mathematical model of infrequent changes of discrete states over time, in which future events occur by chance and depend only on the current state, and not on the history of how that state was reached."
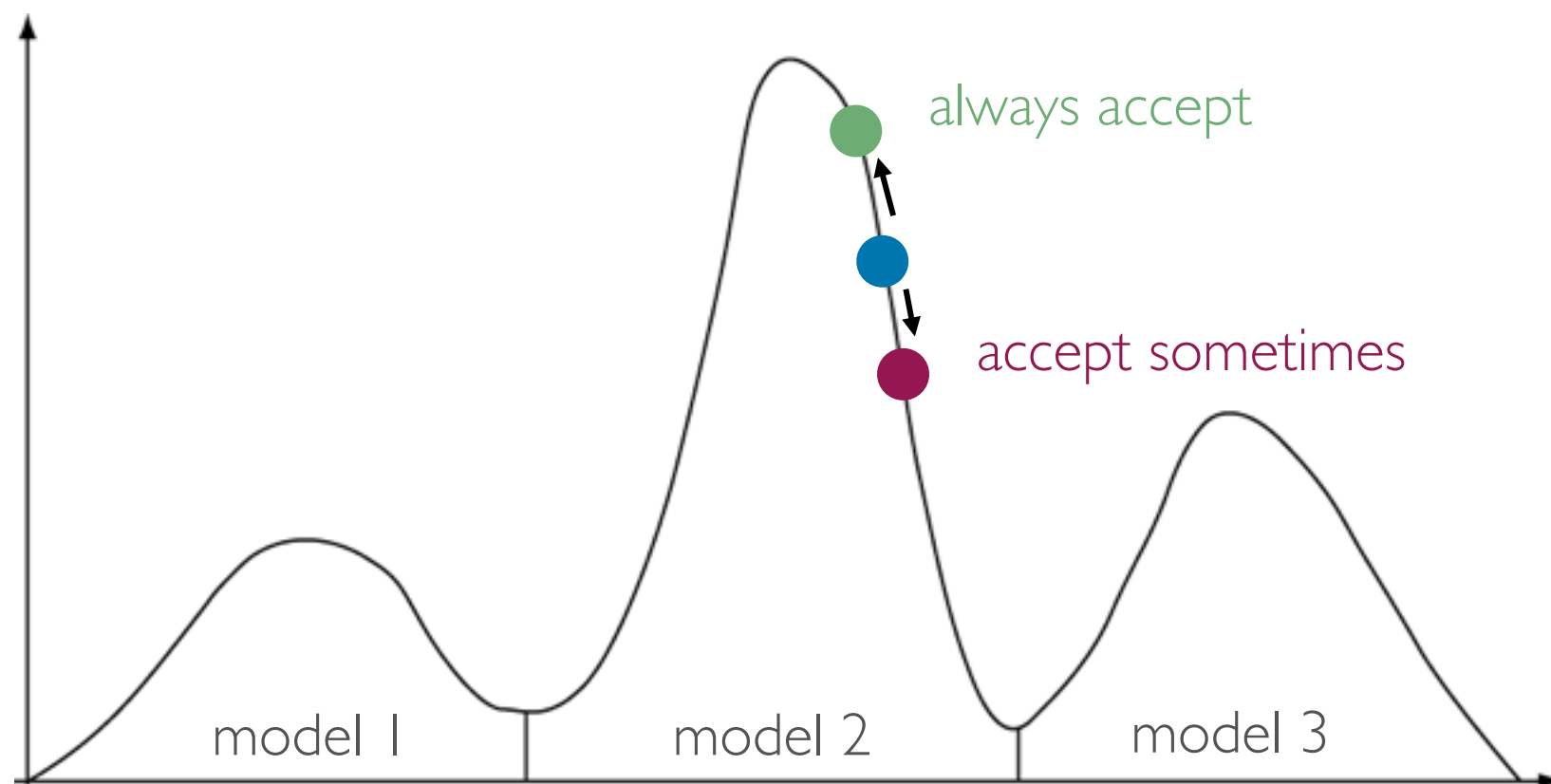
state at $t_{t+1}$ only depends on state $t$

Markov chain
process has a finite state-space
(i.e. you can count it)

Whelan et al 2001

[1] Start at an arbitrary point in parameter space

[2] Make a small random move in one parameter

[3] Calculate proposal ratio (R) of new state to old state:

    [a] if R > 1 then the new state is accepted

    [b] if R < 1 then the new state is accepted with probability R

    [c] If new state not accepted, stay in the old state

always accept

accept sometimes

model 1    model 2    model 3

slide adapted from Frederick Ronquist

[3] Calculate proposal ratio (R) of new state to old state

$$R = \min[1, \text{Lh ratio} \times \text{prior ratio} \times \text{proposal ratio}]$$

[a] if R > 1 then the new state is accepted

[b] Generate random variable U[0,1]*

If R > U then the new state is accepted
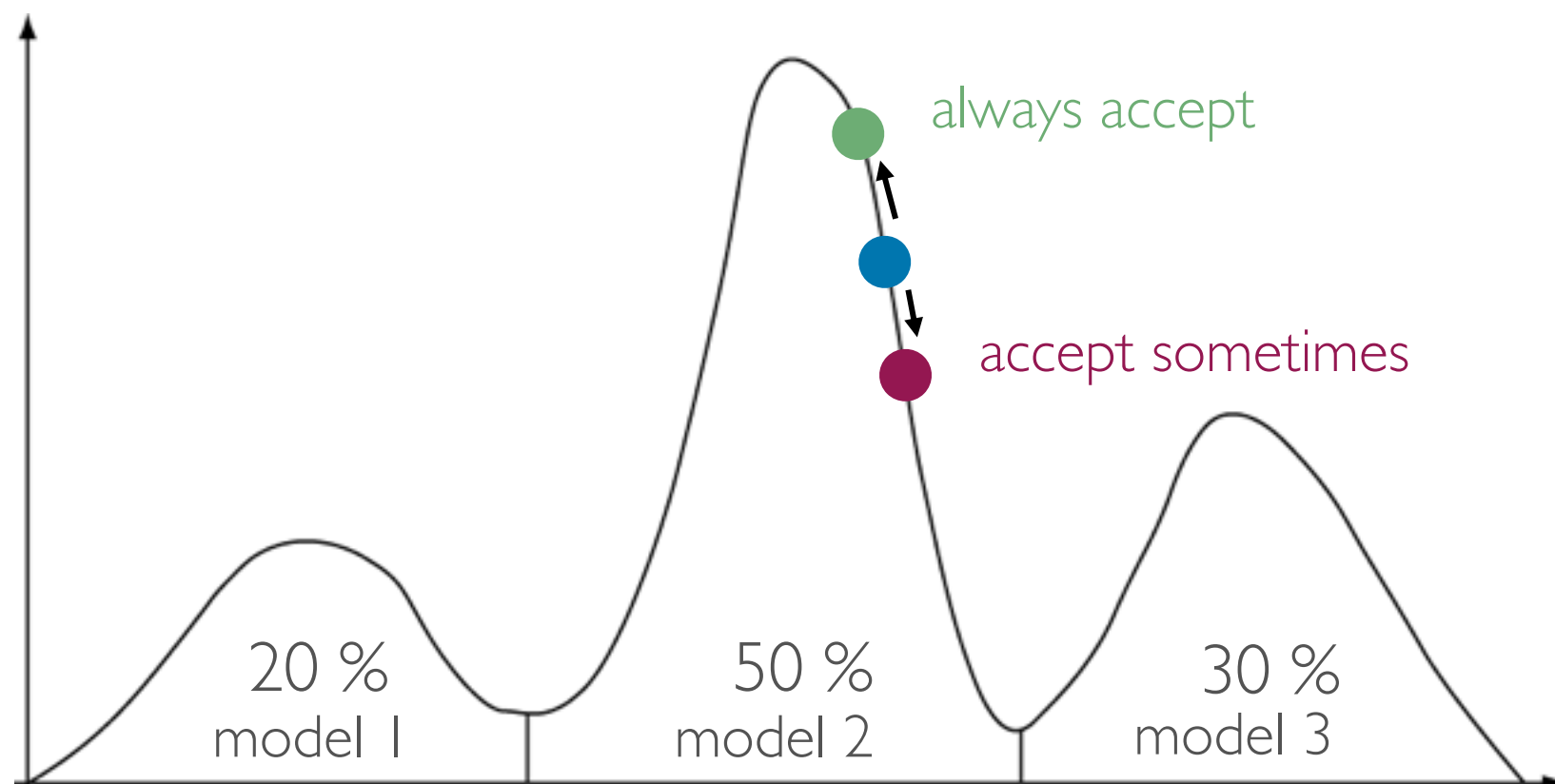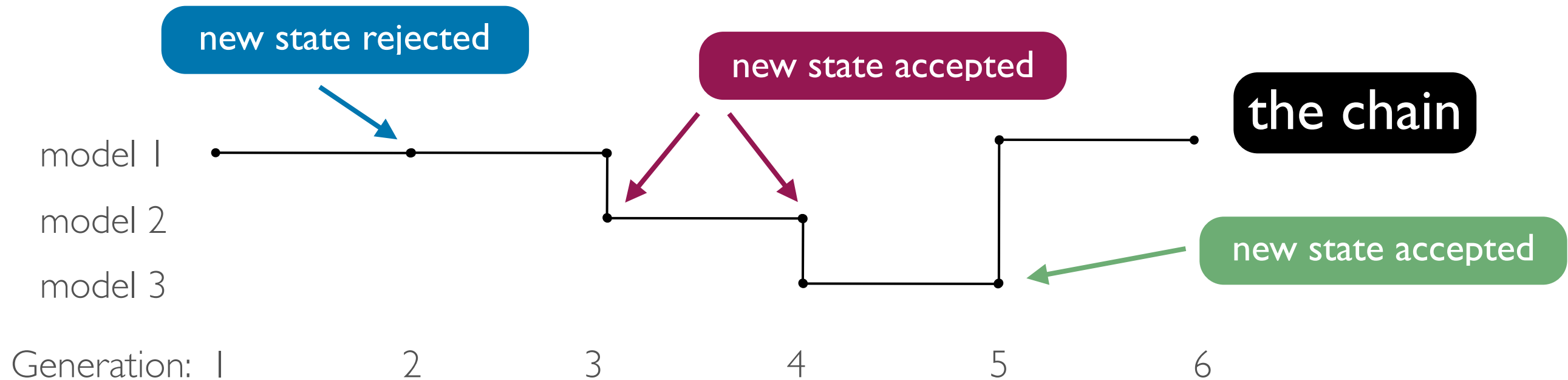
[c] If new state not accepted, stay in the old state

always accept

accept sometimes

20 %
model 1

50 %
model 2

30 %
model 3

* allows us to fully characterise the marginal distribution

slide adapted from Frederick Ronquist

[1] Start at an arbitrary point in parameter space

[2] Make a small random move in one parameter

[3] Calculate proposal ratio (R) of new state to old state:

   Accept or not based on R, U

[4] Repeat many, many times: a Markov chain



always accept

accept sometimes

The proportion of time the MCMC procedure samples from a particular parameter region is an estimate of that region's posterior probability density
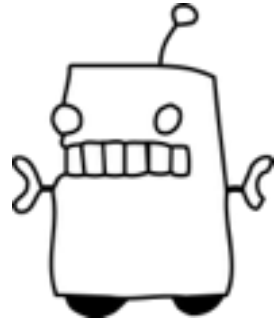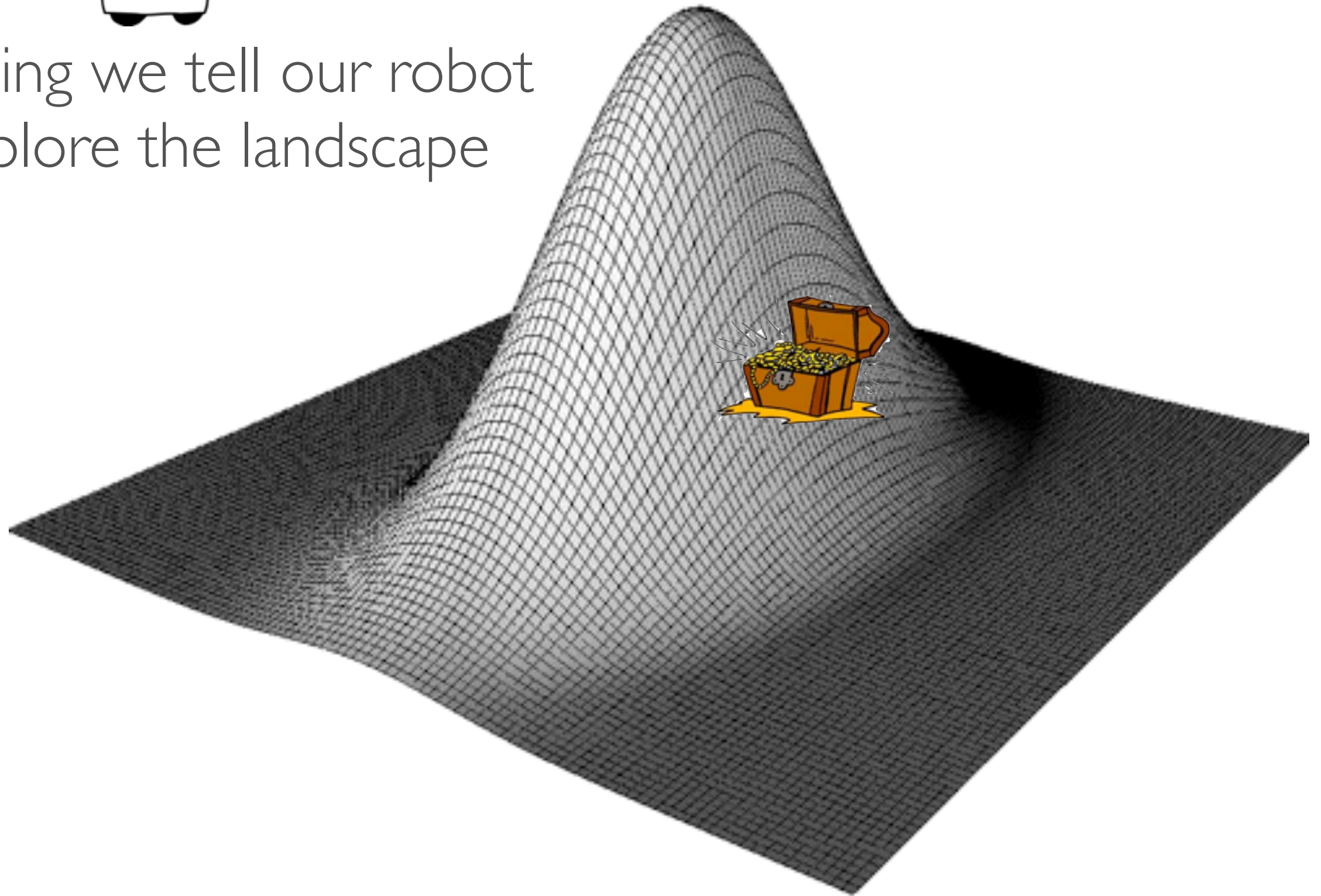
20 %
model 1

50 %
model 2

30 %
model 3

slide adapted from Frederick Ronquist

# Markov chain Monte Carlo sampler 4

new state rejected

new state accepted

the chain

new state accepted

model 1

model 2

model 3

Generation: 1    2    3    4    5    6

Bayesian Posterior Probability for model 1 (BPP$_{model\ 1}$)= 4/6

Sampling the MCMC provides a valid approximation
for the posterior distribution of trees
(over 100,000s – 1,000,000s of generations)
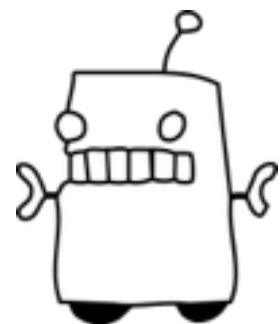without having to know the denominator
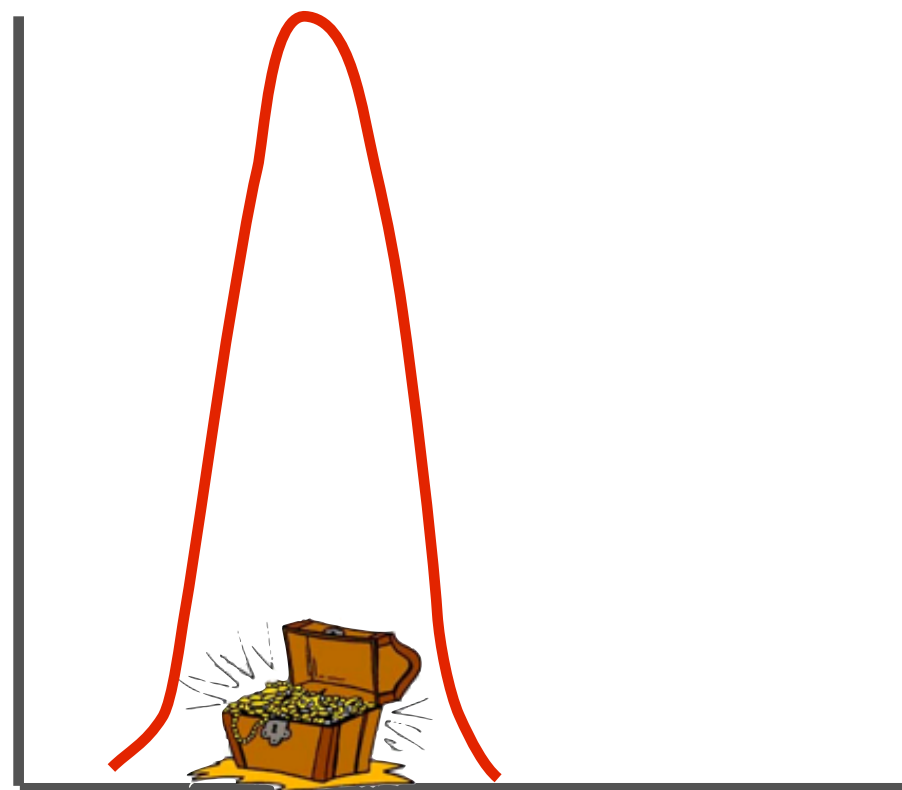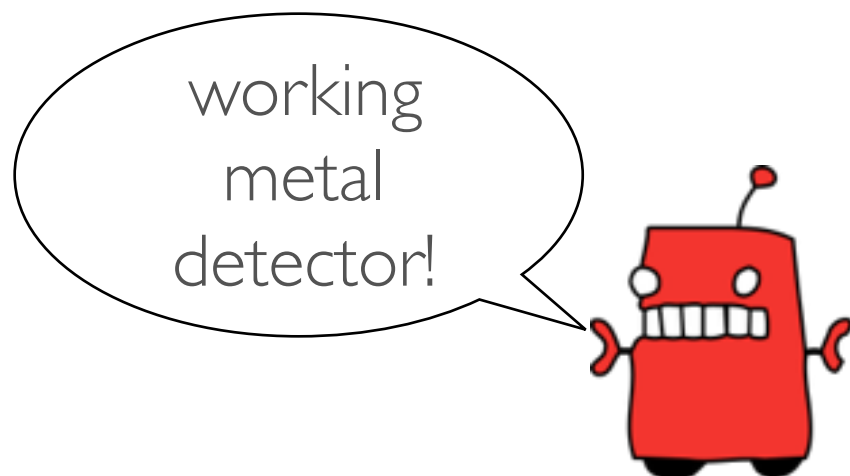
something we tell our robot
to explore the landscape

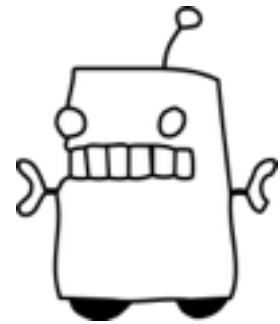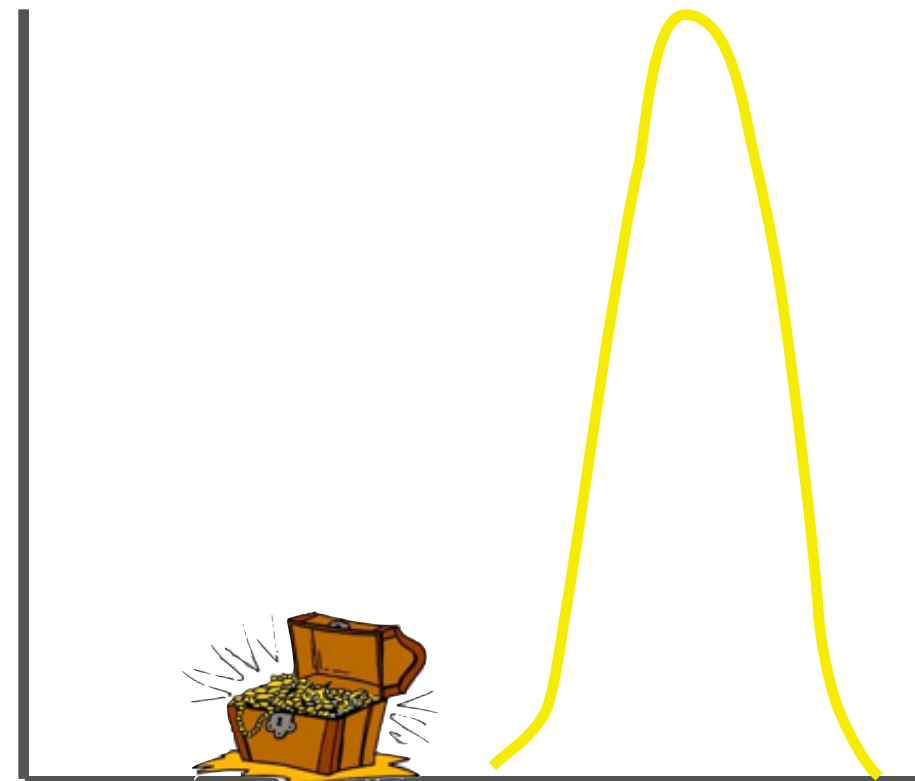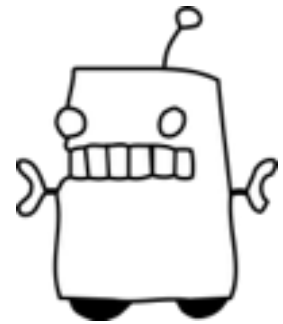information formalised as a distribution

# informative prior

we have some bad information

working glass detector!
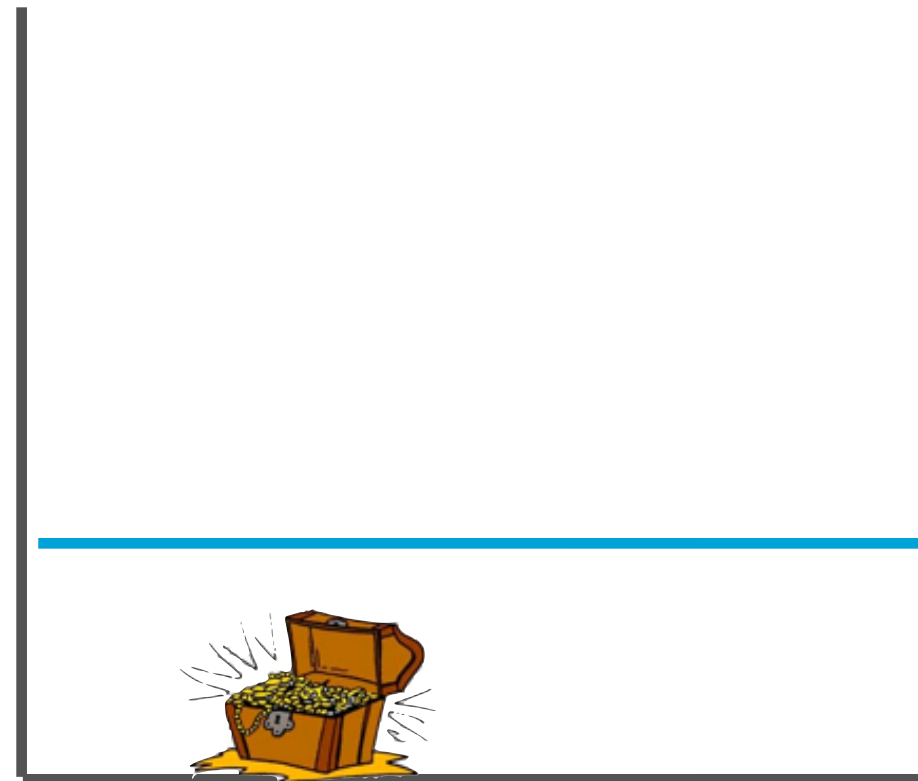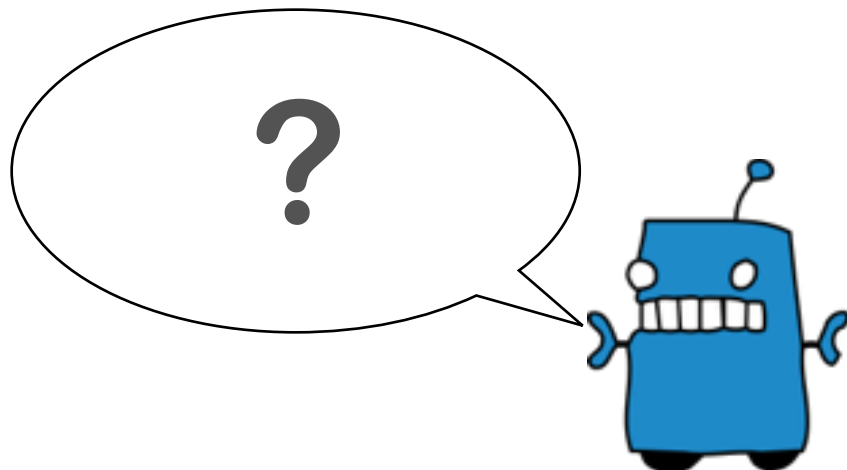
information formalised as a distribution
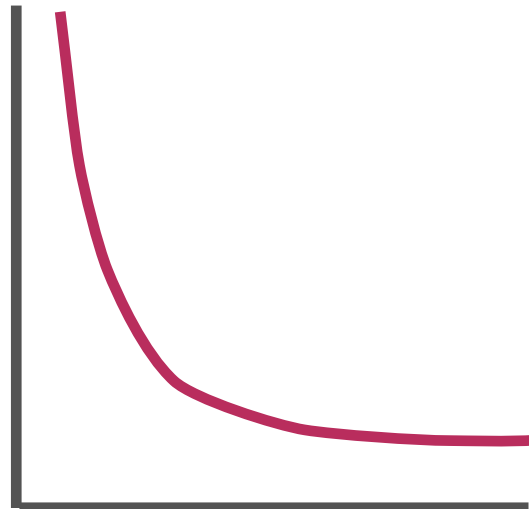
uniform prior

sometimes we don't know anything

?

uniform

exponential

gamma

low **α**     high **α**

Parameters can vary:

(1) across sites (characters)

(2) across the tree (heterotachy)

$q01 = \alpha$
$q10 = \beta$

# rate variation across sites

## residence

$q01 = \alpha = 1$
$q10 = \beta = 2$

$q01 = \alpha = 1$
$q10 = \beta = 2$

## descent

$q01 = \alpha = 1$
$q10 = \beta = 2$

$q01 = \alpha = 4$
$q10 = \beta = 1$

characters have same rates

characters have different rates

## The Gamma Solution:
## a relative rate multiplier for branches

- shrinks branches for faster rates of evolution, stretches for lower
- draw multipliers from a gamma distribution, magically approximated by actually just four rates

Yang 1994

rate variation across the tree

# rate variation across the tree

### residence

$q01 = \alpha = 1$
$q10 = \beta = 2$

$q01 = \alpha = 4$
$q10 = \beta = 1$

### descent

$q01 = \alpha = 1$
$q10 = \beta = 2$

$q01 = \alpha = 4$
$q10 = \beta = 1$

characters have same rates ...
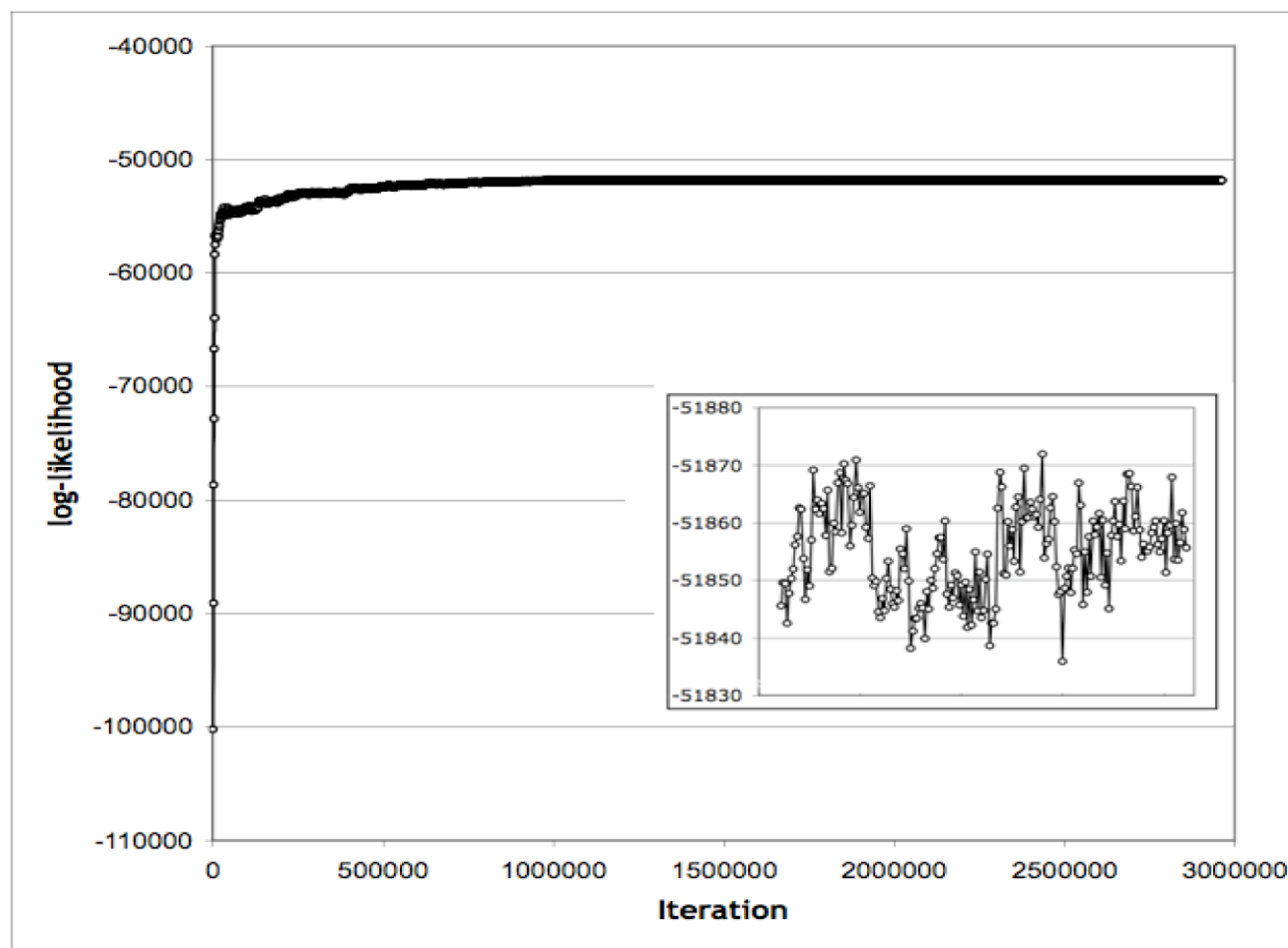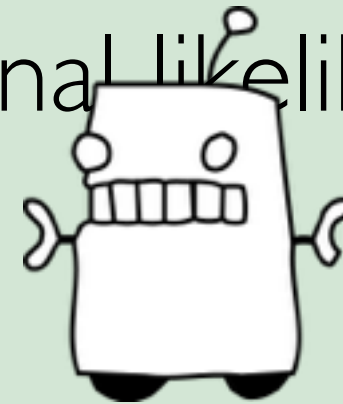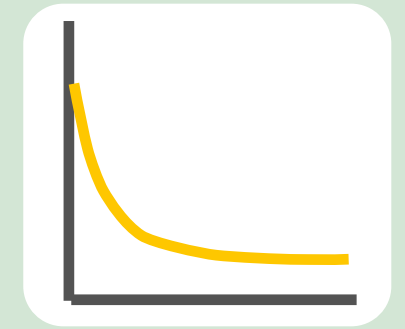
... but rates differ across the tree

## The Covarion Solution:
sites are switched on/off at different parts of the tree

posterior probability distribution $=$ $\dfrac{\text{likelihood of data} \times}{\text{marginal likelihood}}$

# posterior probability distribution 2

## diagnostics

chains
burn-in
convergence
acceptance rates
sampling / ESS
prior influence
$MC^3$

## summary

Lh tracer plots
all parameters
histograms
credibility interval
marginal likelihood

## interpret

Bayes factors
posterior:prior
model plots

**Tracer** is very useful
http://tree.bio.ed.ac.uk/software/tracer/

Tracer: Rambaut & Drummond

[1] As many as you can

[2] As long as you can run it

[3] Sample to reduce autocorrelation

[4] Use $MC^3$ for hot/cold chains in one run if possible

# diagnostics 2: trace plots

**ACCEPTANCE**

[1]  Acceptance of the proposed move has a rate
[2]  Aim for 20-70% acceptance
[3]  Use the tuning parameter to hit this target
      (e.g. *ratedev* in BayesTraits)

**ESS**

[1]  ESS: effective sample size
[2]  Each chain step is correlated with the previous
[3]  Sample from the chain at large-ish intervals to
      reduce autocorrelation
[4]  Plot $\theta_t$ by $\theta_{t+1}$ to assess

[1] As many as you can
[2] As long as you can run it
[3] Sample to reduce autocorrelation
[4] Use MC$^3$ for hot/cold chains in one run if possible

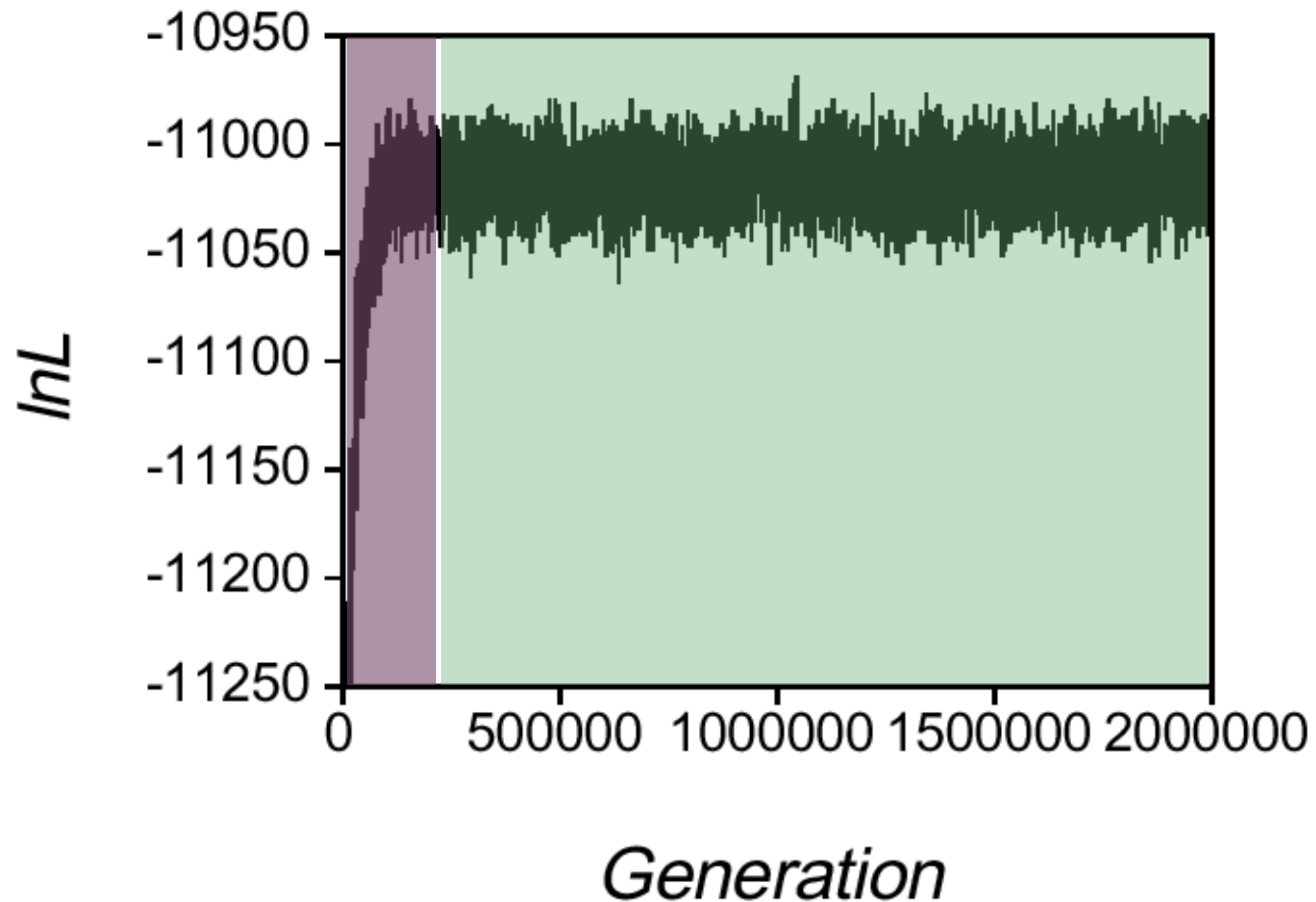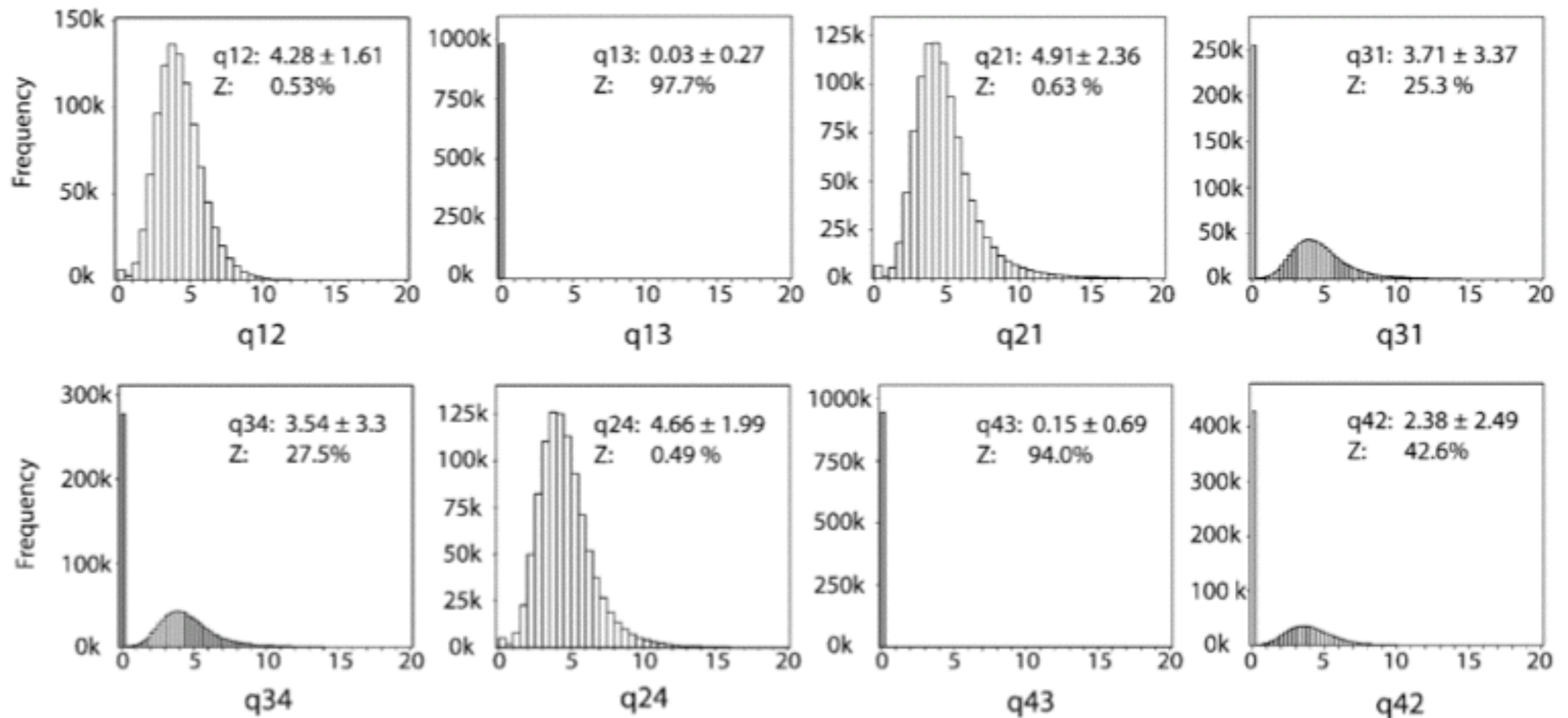95% credibility interval

contains the true values with 95% probability

Approximate the marginal likelihood of the PPD with the **harmonic mean** of the likelihoods

$$H = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \cdots + \frac{1}{x_n}} = \frac{n}{\sum_{i=1}^{n} 1/x_i}, \qquad x_i > 0 \text{ for all } i.$$

**harmonic mean =** reciprocal of the arithmetic mean of all reciprocals

[1] gives more weight to small values
[2] minimises the effect of large values
[3] not without contention!

Kass & Raftery 1995 / Newton & Raftery 1994

$$BF_{ij} = \frac{P(D|M_i)}{P(D|M_j)}$$

$$2\log_e(BF)$$

It can be useful to consider twice the natural logarithm of the Bayes factor, which is on the same scale as the familiar deviance and likelihood ratio test statistics. Rounding and using 20 rather than 10 as the requirement for strong evidence, we then obtain a slight modification:

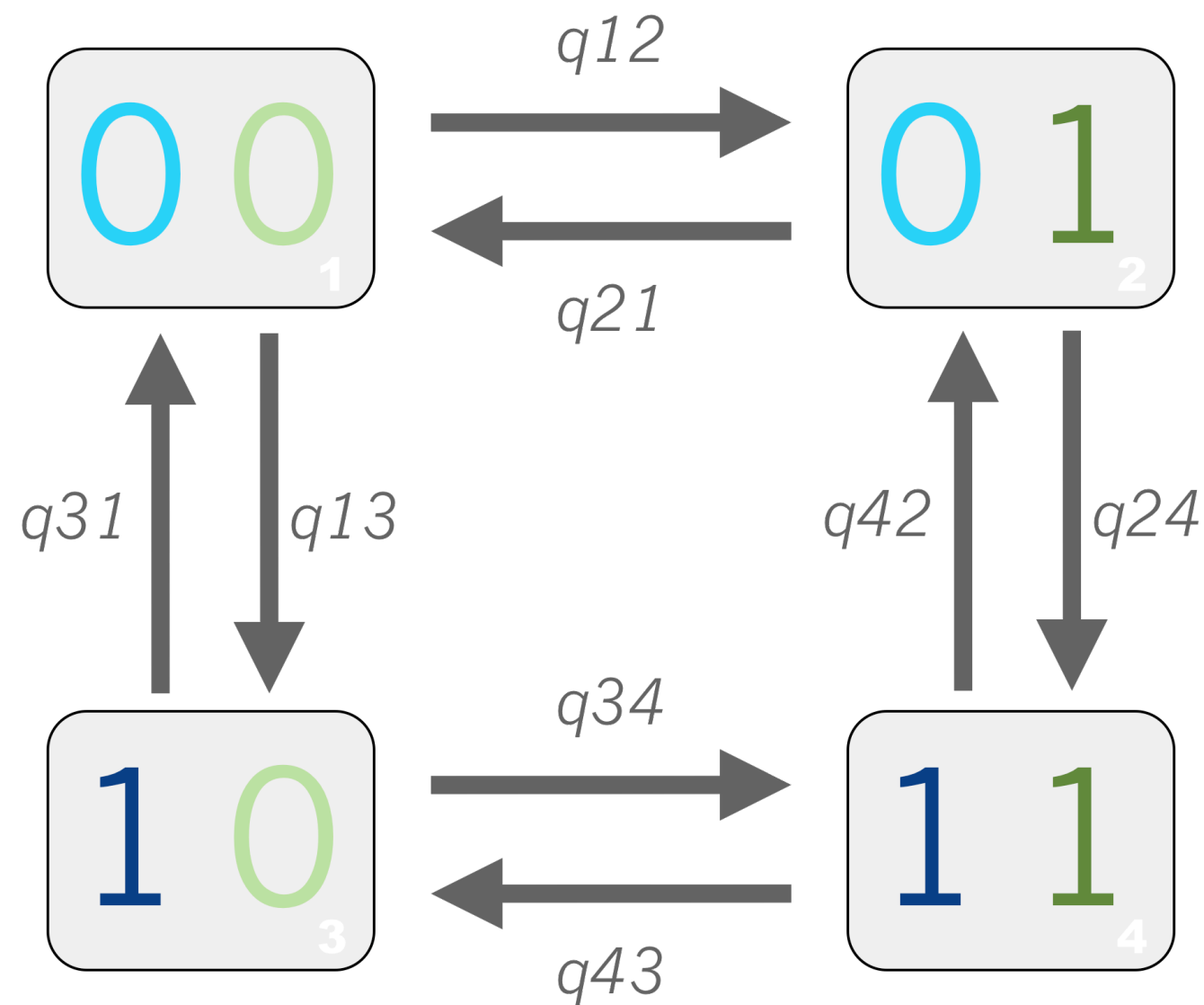| $2 \log_e(B_{10})$ | $(B_{10})$ | Evidence against $H_0$ |
|---|---|---|
| 0 to 2 | 1 to 3 | Not worth more than a bare mention |
| 2 to 6 | 3 to 20 | Positive |
| 6 to 10 | 20 to 150 | Strong |
| >10 | >150 | Very strong |

From our own experience, these categories seem to furnish appropriate guidelines.

# Bayes factor

## An evaluation of the support for one model over another.

No penalties are needed for extra parameters, because a more parameter-rich model has a larger parameter space and therefore a lower prior probability density.
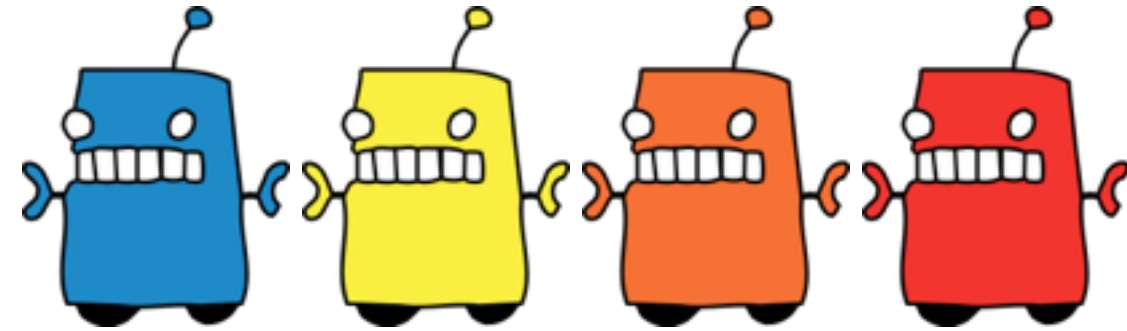
Kass & Raftery 1995 / Newton & Raftery 1994 / Raftery 1996

# Metropolis-coupled MCMC (MC3)

Imagine the PPD surface is made of wax

Multiple MCMC robots descend

A "hot" chain robot acts as a scout for the "cold" captain robot

Heat distorts the joint PPD and flattens it
> this makes it easier to explore regions of low probability

The cold chain (which could be stuck in a local optimum) can escape when a proposed swap with a hot chain is successful.