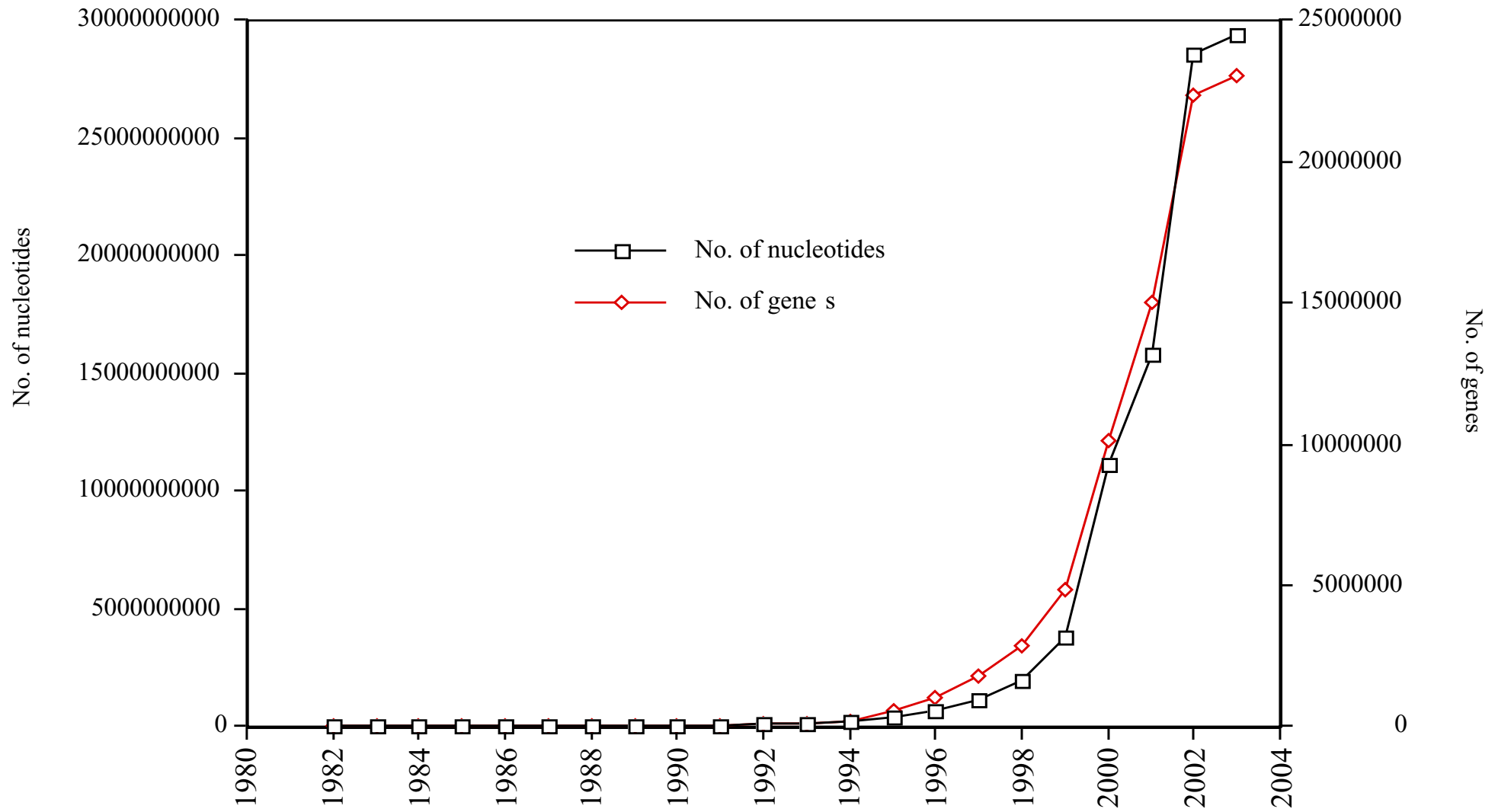# A Model of Pattern-Heterogeneity for Inferring Phylogenetic Trees and Investigating Sequence Evolution

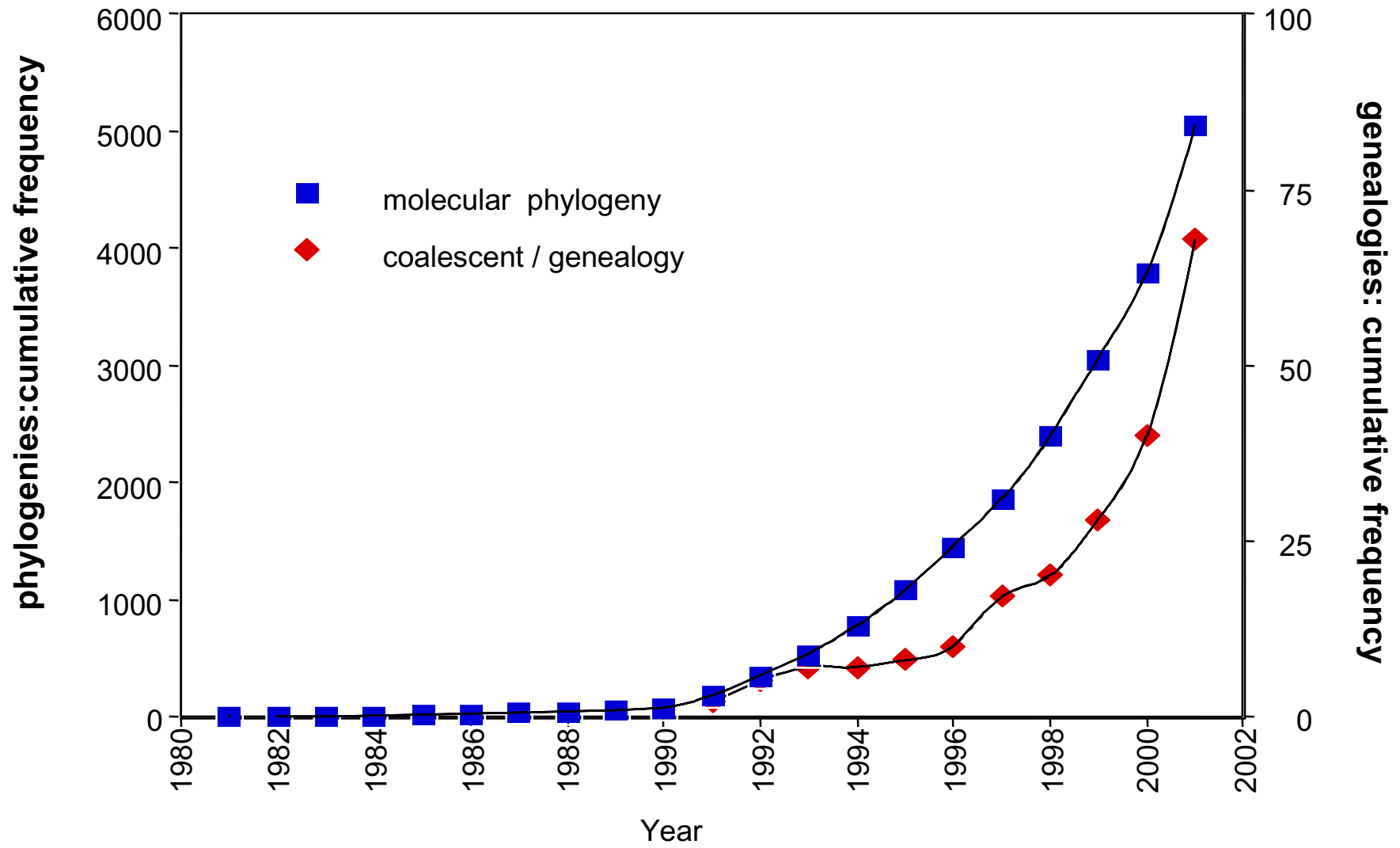## Mark Pagel and Andrew Meade
## Reading University

**m.pagel@rdg.ac.uk**

The growth of GenBank
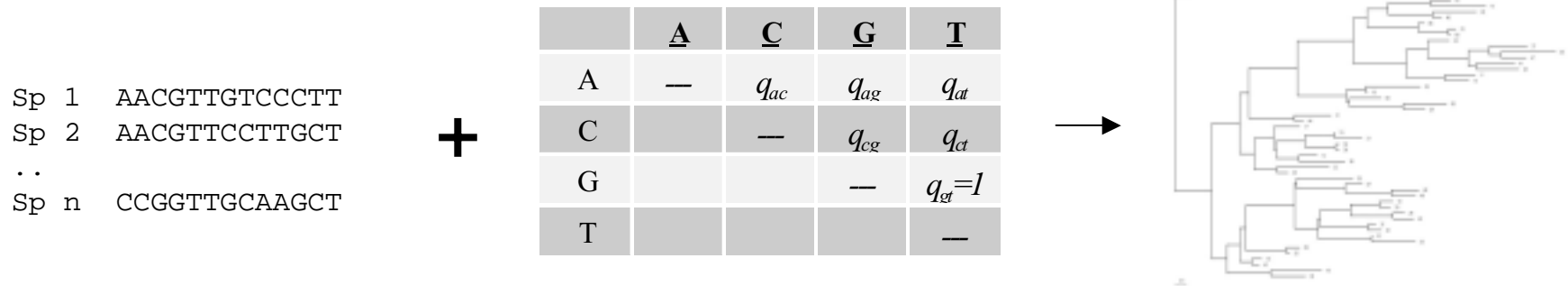
The use of molecular phylogenies in biological research

# Overview:  Phylogenetic Inference from gene sequences

## 1.  Homogeneous Process Model:

Sp 1   AACGTTGTCCCTT
Sp 2   AACGTTCCTTGCT
..
Sp n   CCGGTTGCAAGCT

**+**

|   | **A** | **C** | **G** | **T** |
|---|---|---|---|---|
| A | — | $q_{ac}$ | $q_{ag}$ | $q_{at}$ |
| C |   | — | $q_{cg}$ | $q_{ct}$ |
| G |   |   | — | $q_{gt}=1$ |
| T |   |   |   | — |



## 2. Modifications to basic homogeneous model

a) **rate-heterogeneity** (Yang, 1994):  apply homogeneous model but allow rates to vary over sites according to a discrete gamma.  Equivalent to fitting $k$ rate matrices to each site, where the matrices differ from each other only by a proportional scalar

b) **partitioning data**:  apply a different rate matrix to different sites, chosen beforehand by the investigator

**Limitations to homogeneous, rate-heterogeneity and partitioning approaches**

Homogeneous model:  not all sites may evolve according to the same model

Gamma model:  rate variation may not be 'gamma' or sites may vary in some other way

Partitioning data: presumes investigator knows with certainty which model best applies to each site

# Pattern-Heterogeneity Model of Gene-Sequence Evolution

Allow for different sites to evolve in *qualitatively* (or quantitatively) different ways without prior partitioning by the investigator.

Method: fit more than one model of evolution to each site, summing the likelihood over all models.  Allow the data to determine the 'best' model for each site.

Motivation for model:  *Pattern-heterogeneity* model will always equal or better the performance of homogeneous, gamma rate heterogeneity and partitioning models.  Frequently yields substantial improvements (100's of log-units)


Applications

Phylogenetic inference

Detecting regions within genes that evolve differently

Detecting differences among genes

# An example of pattern-heterogeneity

## Sequence data

```
Sp 1    AACGTTGTCCCTT
Sp 2    AACGTTCCTTGCT
..
Sp n    CCGGTTGCAAGCT
```
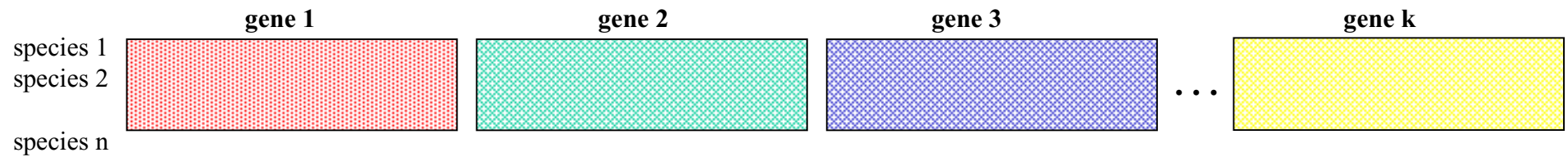
## Two transition rate matrices

|   | **A** | **C** | **G** | **T** |
|---|---|---|---|---|
| A | — | 4.97 | 3.41 | 0.82 |
|   |   | 2.11 | 3.13 | 0.34 |
| C |   | — | 0.35 | 2.82 |
|   |   |   | 3.87 | 1.49 |
| G |   |   | — | 1 |
|   |   |   |   | 1 |
| T |   |   |   | — |

# Applications of pattern-heterogeneity model

**Single gene alignment**

species 1

**pattern 1**  **pattern 2**  **pattern 3**

species n

**Concatenated gene alignment**

**gene 1**  **gene 2**  **gene 3**  **gene k**

species 1
species 2

species n

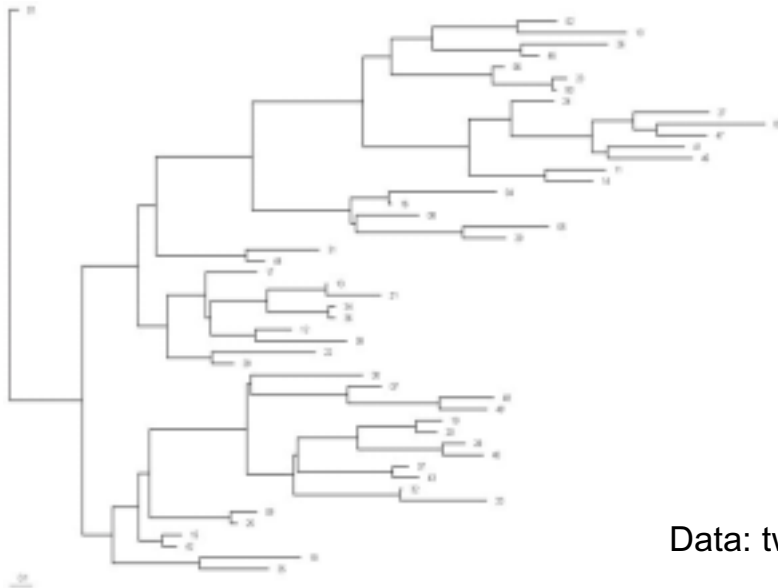**"Supermatrix" alignment**

species 1
species 2

species n-k

species n

# Testing the Pattern-Heterogeneity Model: detecting pattern-heterogeneity in simulated data

Method:

I) generate simulated gene-sequence data on a random tree according to two different models of sequence evolution, creating two 'genes' with different patterns of substitutions

2) analyse simulated data using homogeneous, gamma rates and pattern-heterogeneity model.  Evaluate trees by their likelihood

Tree used in simulations

rate matrices



| | **A** | **C** | **G** | **T** |
|---|---|---|---|---|
| A | — | 4.97 | 3.41 | 0.82 |
| | | 2.11 | 3.13 | 0.34 |
| C | | — | 0.35 | 2.82 |
| | | | 3.87 | 1.49 |
| G | | | — | 1 |
| | | | | 1 |
| T | | | | — |

Data: two 'genes' of length 1200 and 800 sites

Brief digression….

Markov Chain Monte Carlo methods for inferring phylogenies

1 Construct a markov chain whose successive states are possible phylogenetic trees

2. Start at a random tree then guide the chain such that it samples a desirable region of the universe of possible trees (Metropolis-Hastings ratio)

3. Use the sample of trees to estimate parameters of the model of evolution and features of the tree itself

# Sampling the Universe of Phylogenetic Trees

Markov-Chain Monte Carlo (MCMC) Methods

• Generate a large number of phylogenetic trees from a Markov Chain
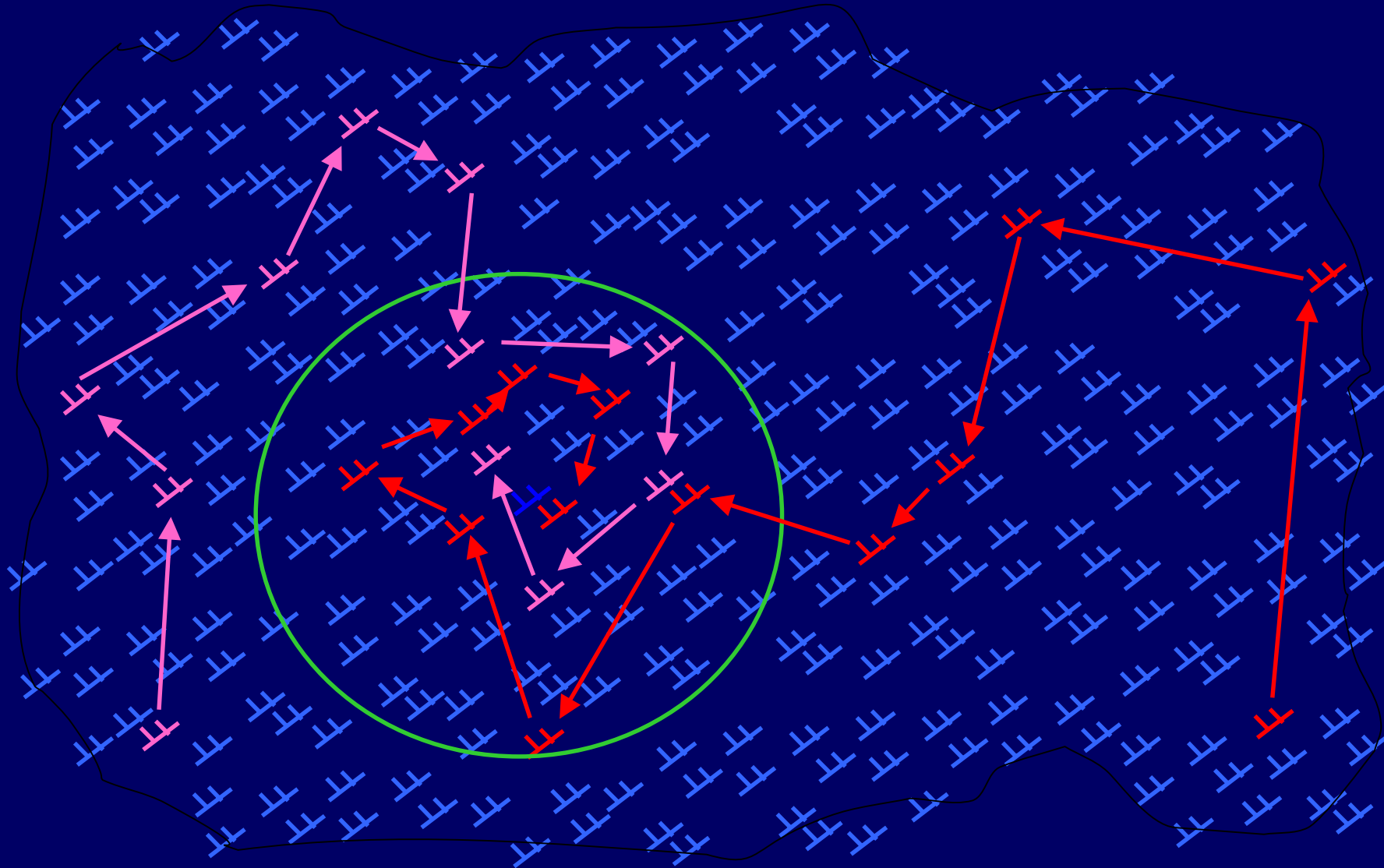• at equilibrium randomly sample from universe of trees

sampling mechanism: The Metropolis-Hastings Algorithm
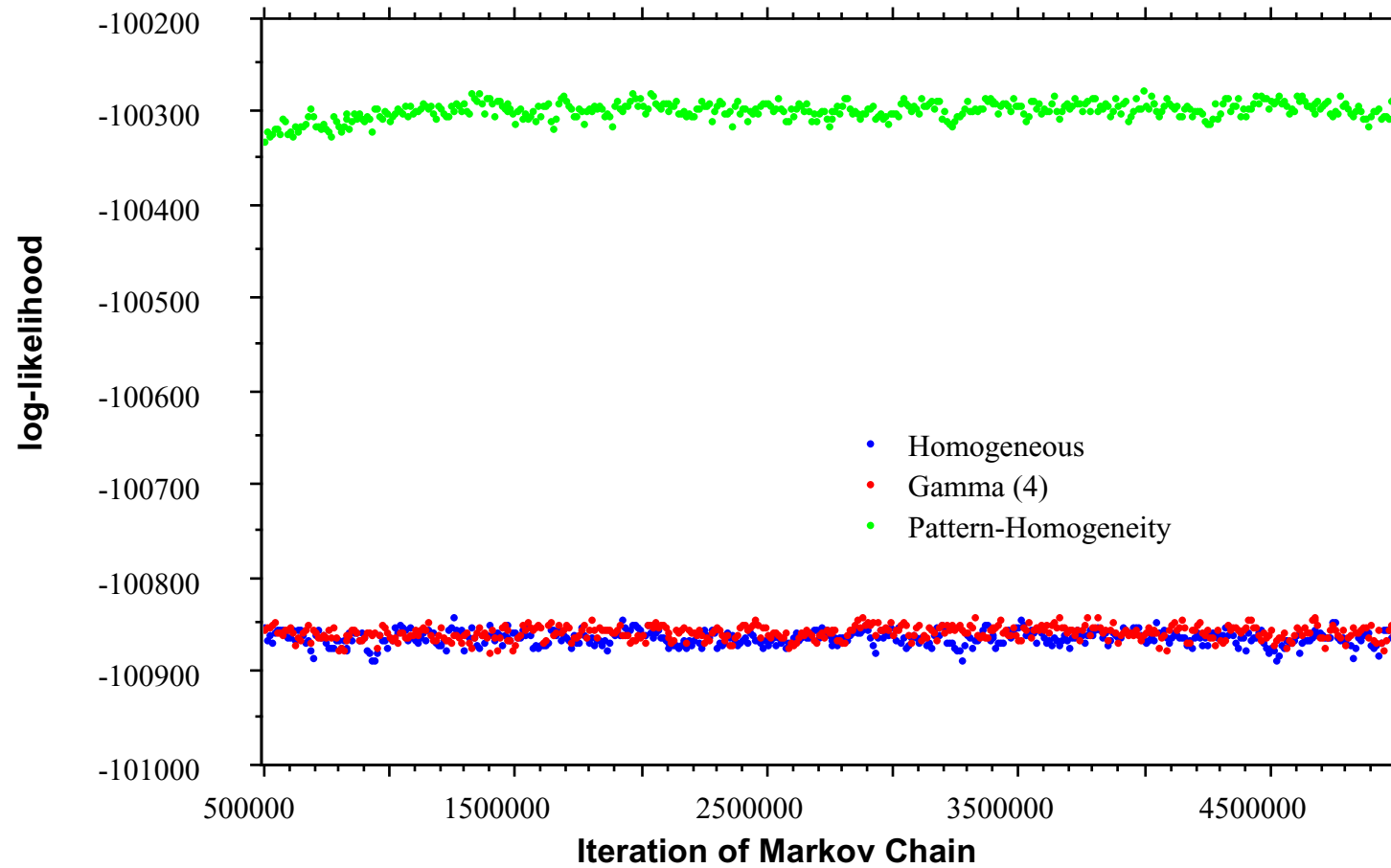
Accept new tree with p=1.0 if $L(T_{n+1}) > L(T_n)$

otherwise…

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

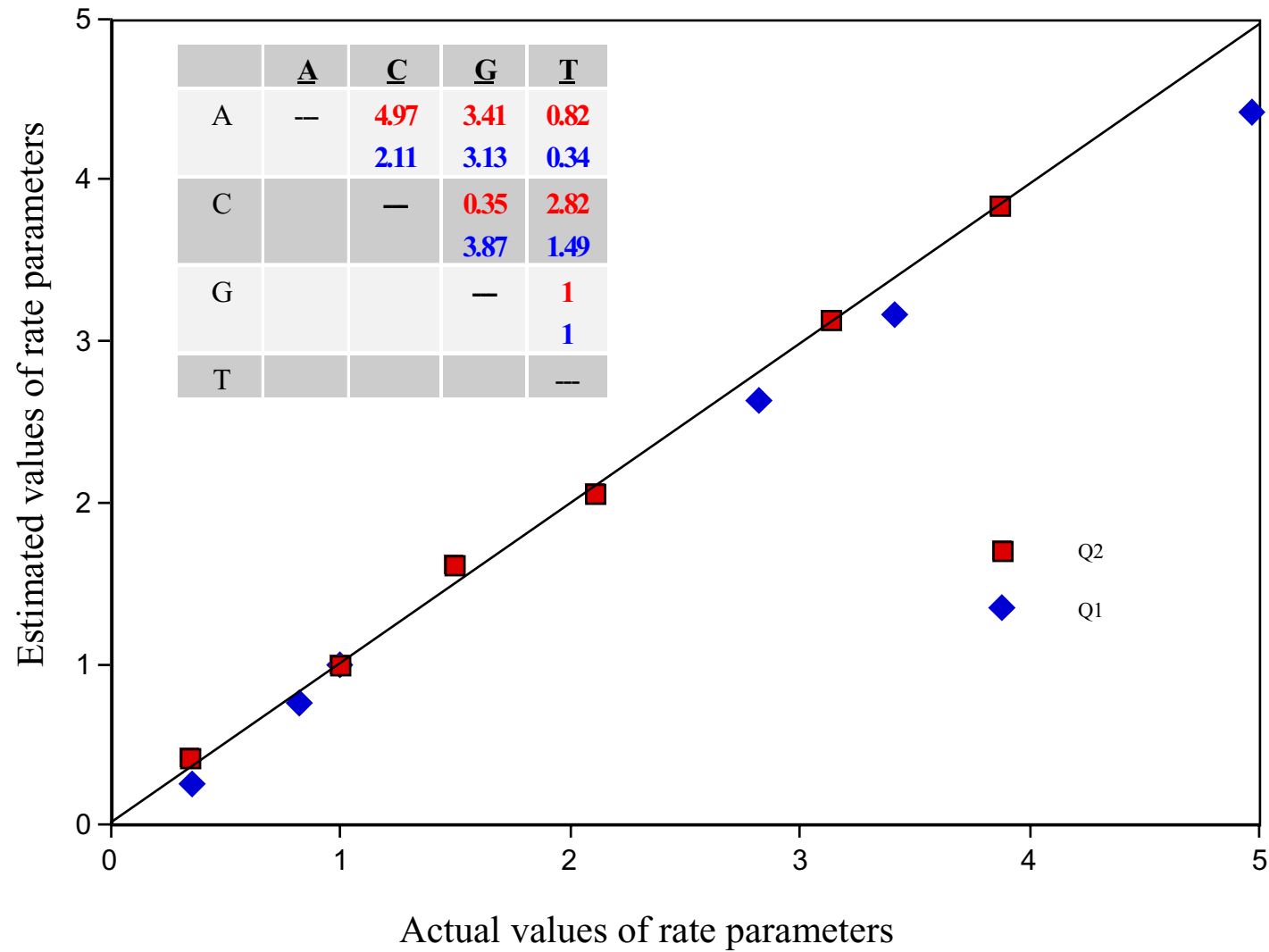Sampling the universe of possible trees:
Markov-chain Monte Carlo methods

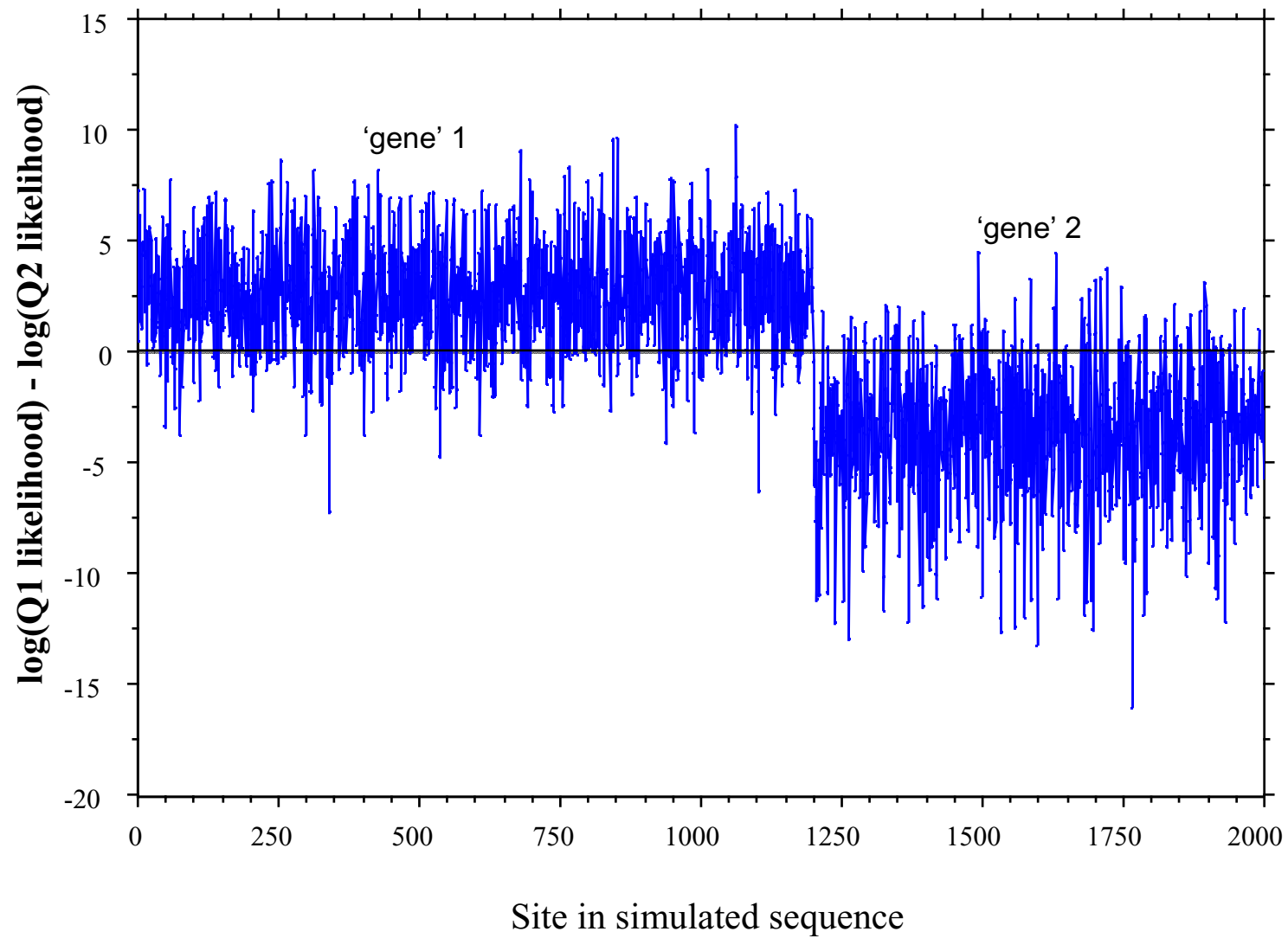**Testing Pattern-Heterogeneity Model: simulations from two random transition rate matrices**

| | A <--> C | A <--> G | A <--> T | C <--> G | C <--> T | G <--> T |
|---|---|---|---|---|---|---|
| Q1 | 4.97 | 3.41 | 0.82 | 0.35 | 2.82 | 1 |
| Q2 | 2.11 | 3.13 | 0.34 | 3.87 | 1.49 | 1 |

**Pattern-heterogeneity model: Simulated and estimated values of the rate parameters**

|   | **A** | **C** | **G** | **T** |
|---|-------|-------|-------|-------|
| A | --- | 4.97 | 3.41 | 0.82 |
|   |   | 2.11 | 3.13 | 0.34 |
| C |   | — | 0.35 | 2.82 |
|   |   |   | 3.87 | 1.49 |
| G |   |   | — | 1 |
|   |   |   |   | 1 |
| T |   |   |   | --- |

Detecting pattern-heterogeneity in simulated data

# Testing the Pattern-Heterogeneity Model: detecting rate-heterogeneity in simulated data

Method:

I) generate simulated gene-sequence data on a random tree according to a gamma rate heterogeneity model (continuous-gamma, α=1.0)

2) analyse data using homogeneous model, gamma rates (4 categories) and pattern-heterogeneity model (4 rate matrices)
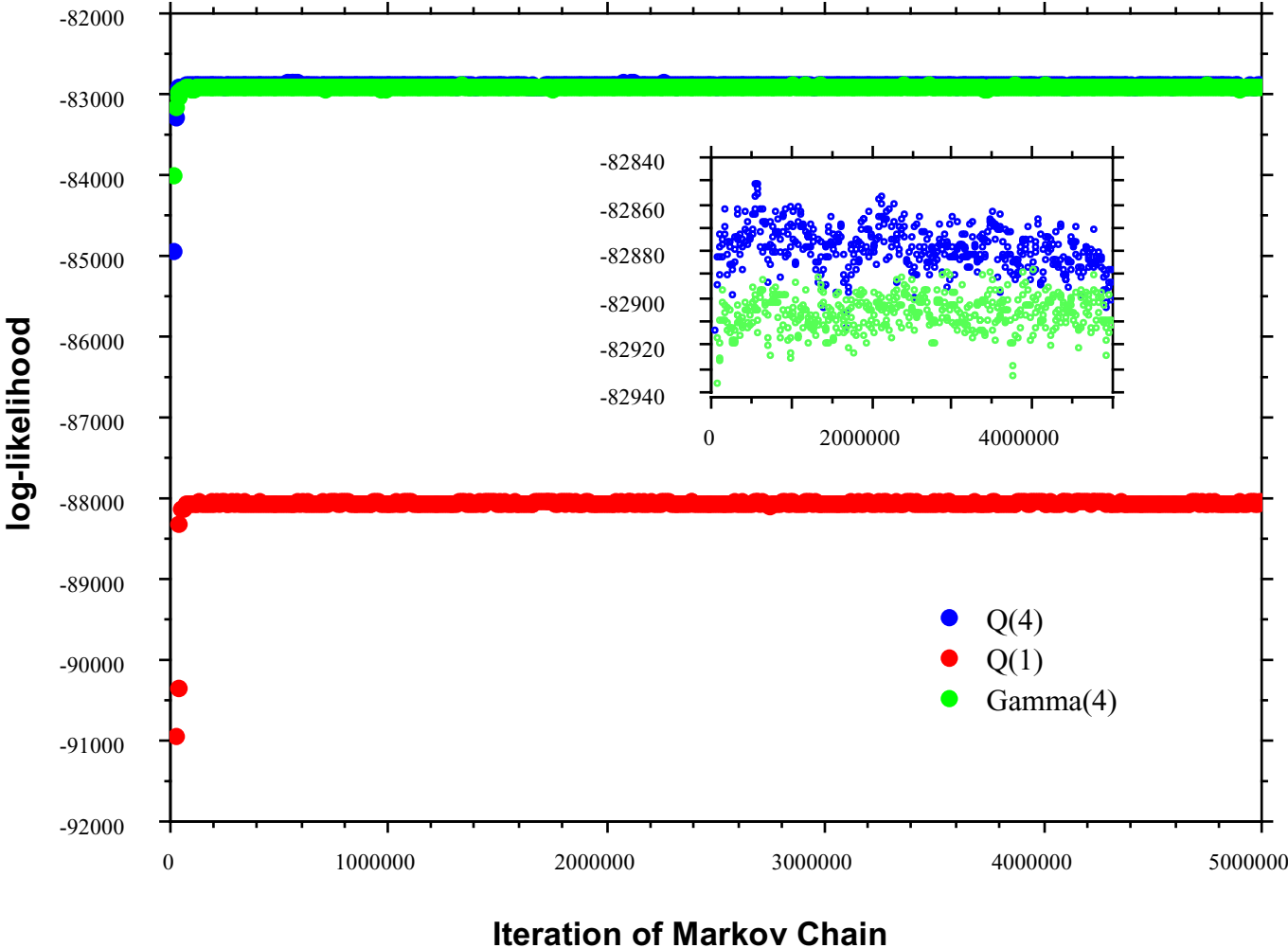
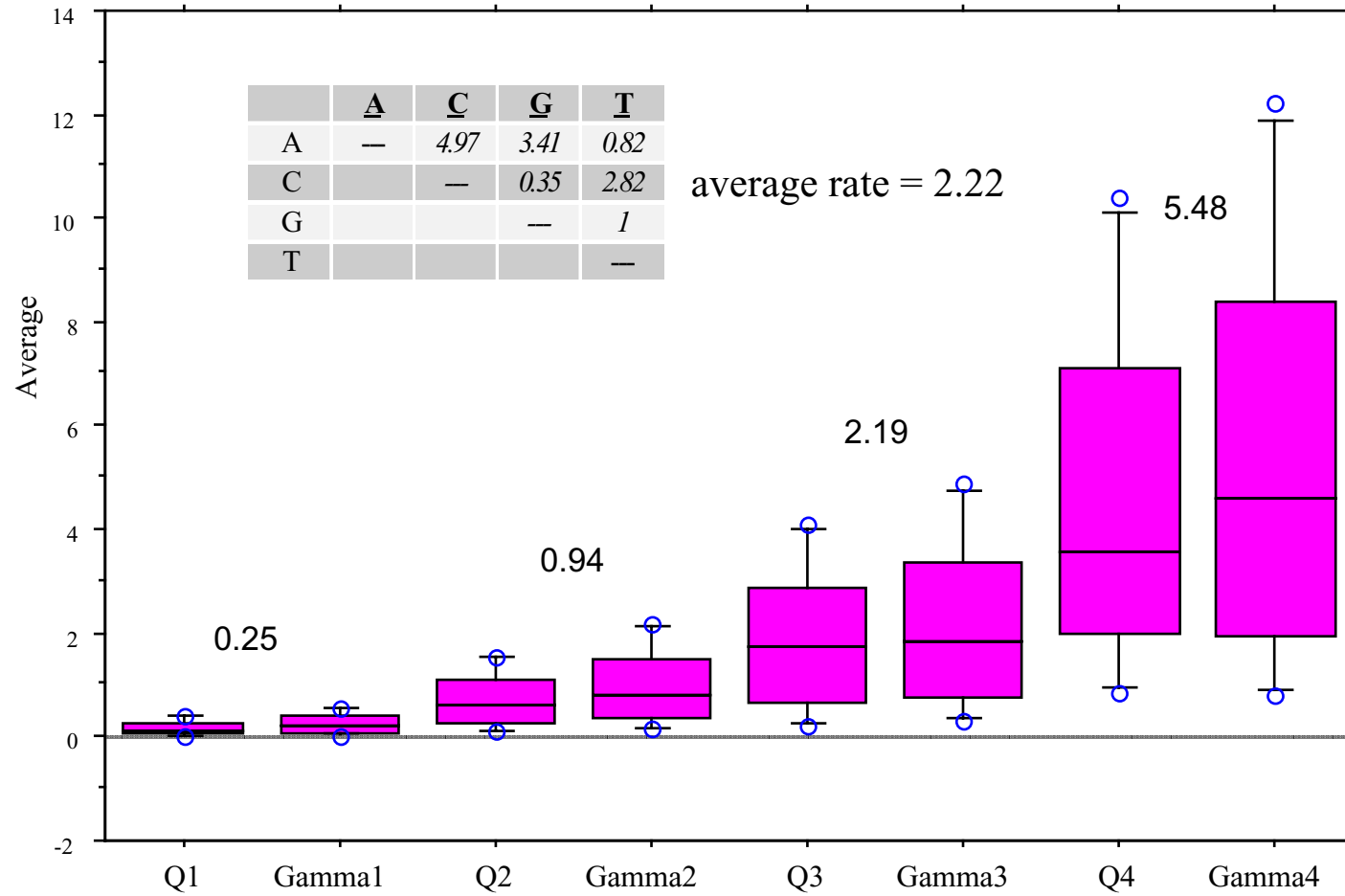Tree used in simulations

rate matrices



|   | **A** | **C** | **G** | **T** |
|---|---|---|---|---|
| A | — | *4.97* | *3.41* | *0.82* |
| C |   | — | *0.35* | *2.82* |
| G |   |   | — | *1* |
| T |   |   |   | — |

Data: a gene of length 2000 sites

**Testing the Pattern-Heterogeneity Model: simulated gamma rate variation**

Average estimated transition rates from gamma and pattern-heterogeneity models applied to simulated gamma rate heterogeneity data: input values shown
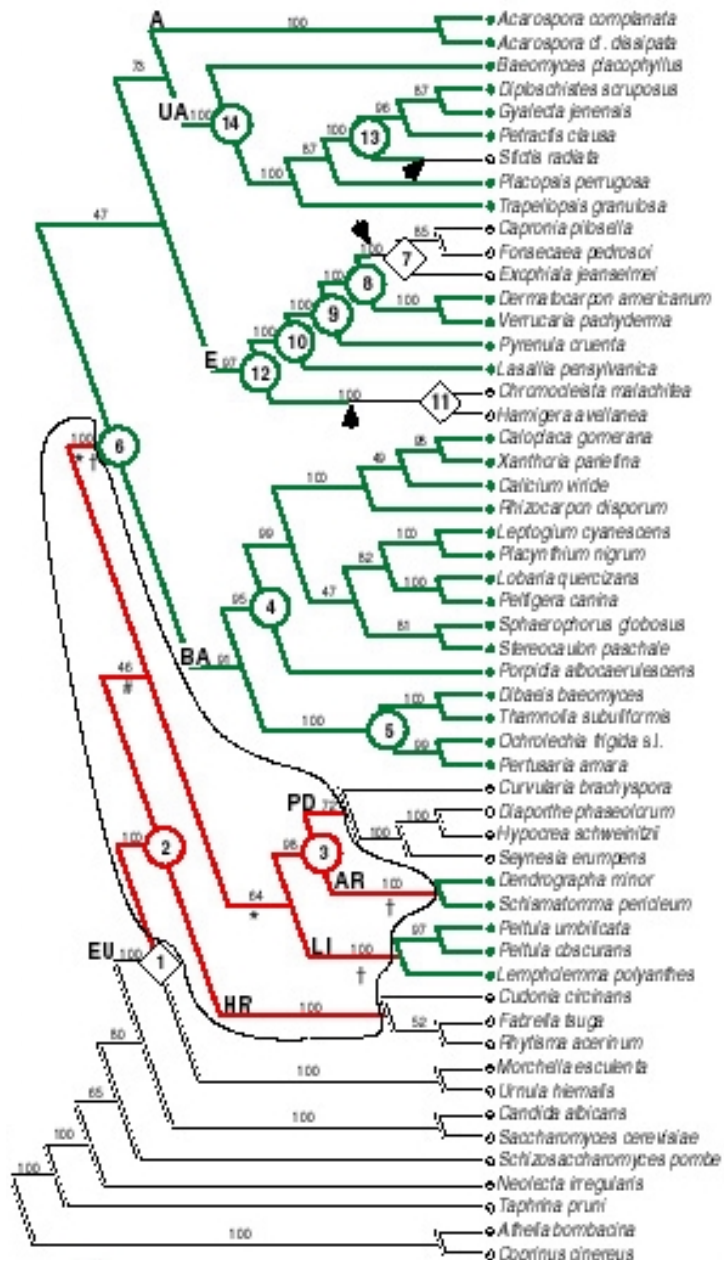
|   | **A** | **C** | **G** | **T** |
|---|---|---|---|---|
| A | — | 4.97 | 3.41 | 0.82 |
| C |   | — | 0.35 | 2.82 |
| G |   |   | — | 1 |
| T |   |   |   | — |

average rate = 2.22

**Applications of the pattern-heterogeneity model to two real data sets**

1. SSU and LSU nrDNA data from Ascomycota fungi: gene differences?


2. Detecting secondary structure in Mitochondrial 12s data from mammals

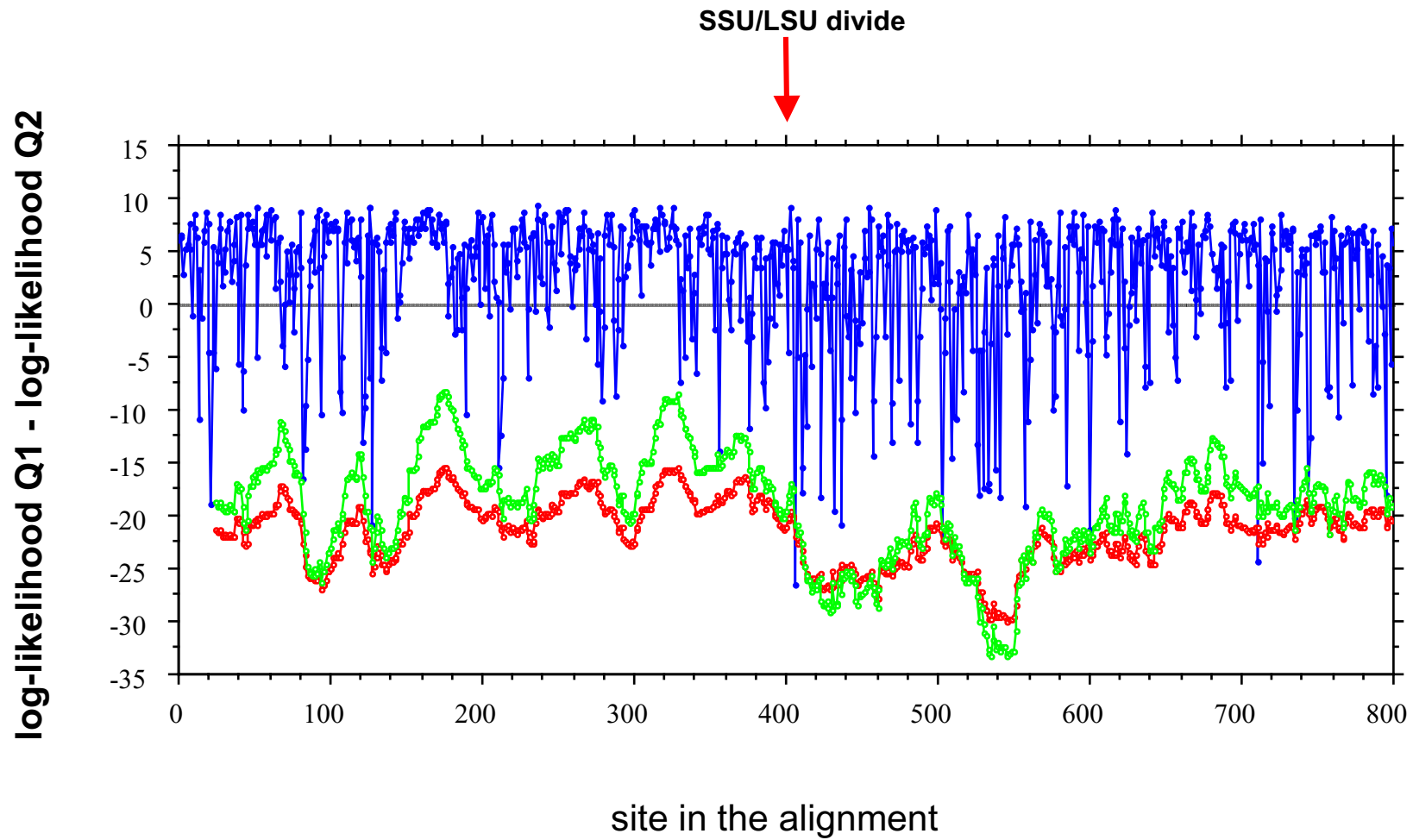**Phylogeny of the Ascomycota Fungi showing the evolution of lichen-formation**

lichen forming

ambiguous

not lichen forming

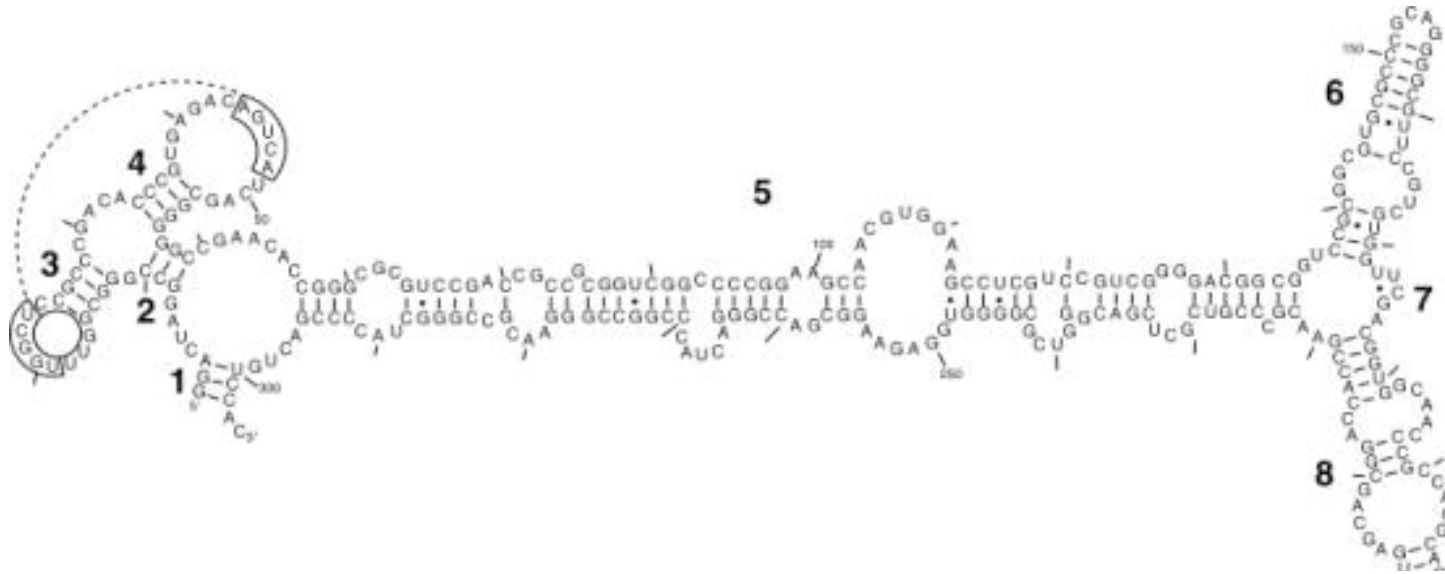log-likelihoods for combined LSU/SSU nrDNA data set:  54 Ascomycota species n=800 sites

# Site by site analysis fitting two independent rate matrices: LSU/SSU nrDNA Ascomycota combined data set



Note: data modelled with two independent rate matrices

# Detecting secondary structure using the pattern-heterogeneity model
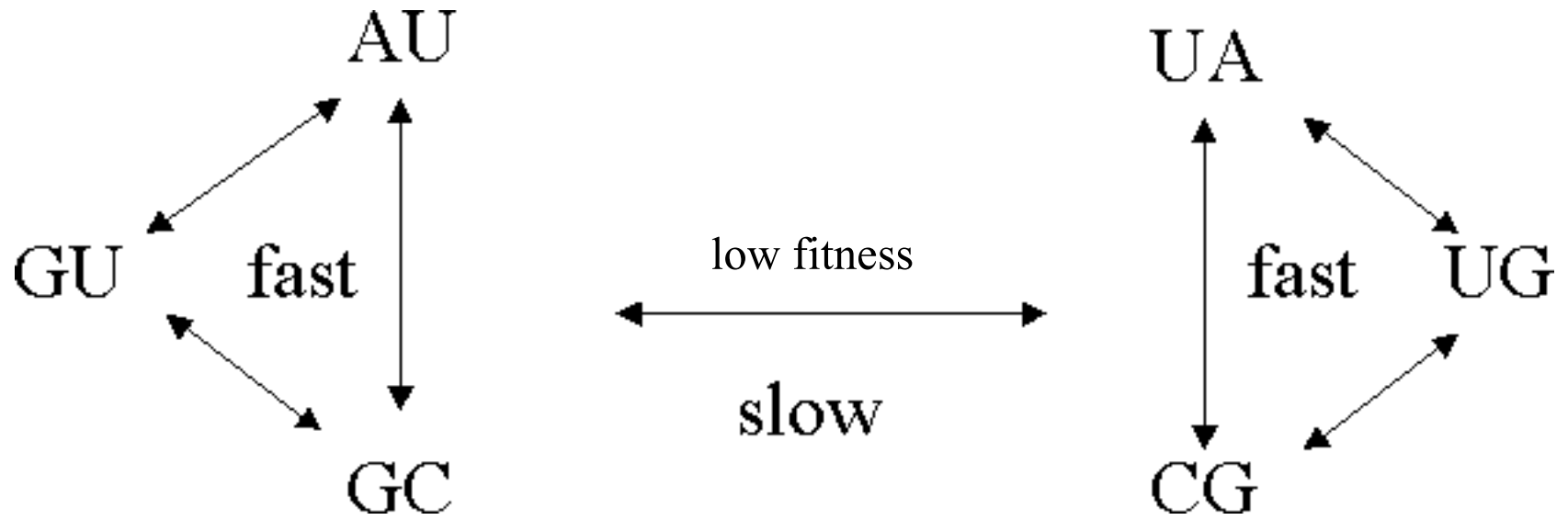


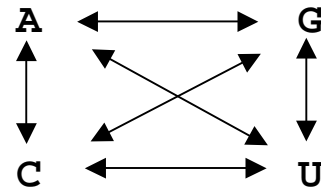Stem and loop structure leads to predictions about the pattern of nucleotide substitutions

   **stems**: dominated by Watson-Crick base pairings.  Therefore expect compensatory substitutions to maintain Watson-Crick pairings.  Non-compensatory changes have lower fitness.  This predicts that transitional changes will occur at a far higher rate than transversional changes.  Often observe Tr/Tv ratio of 10-20 in mitochondrial DNA

   **loops**: no a priori substitutional pattern expected
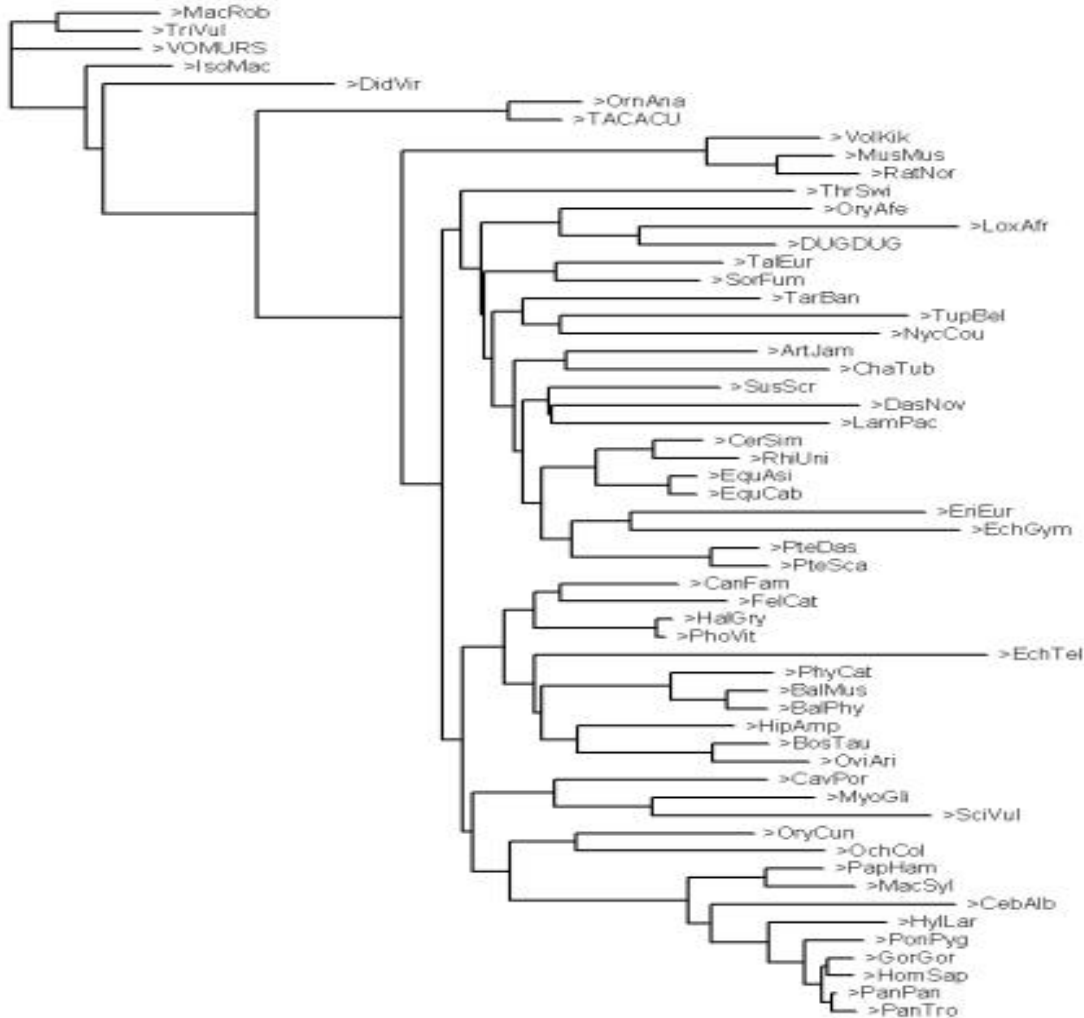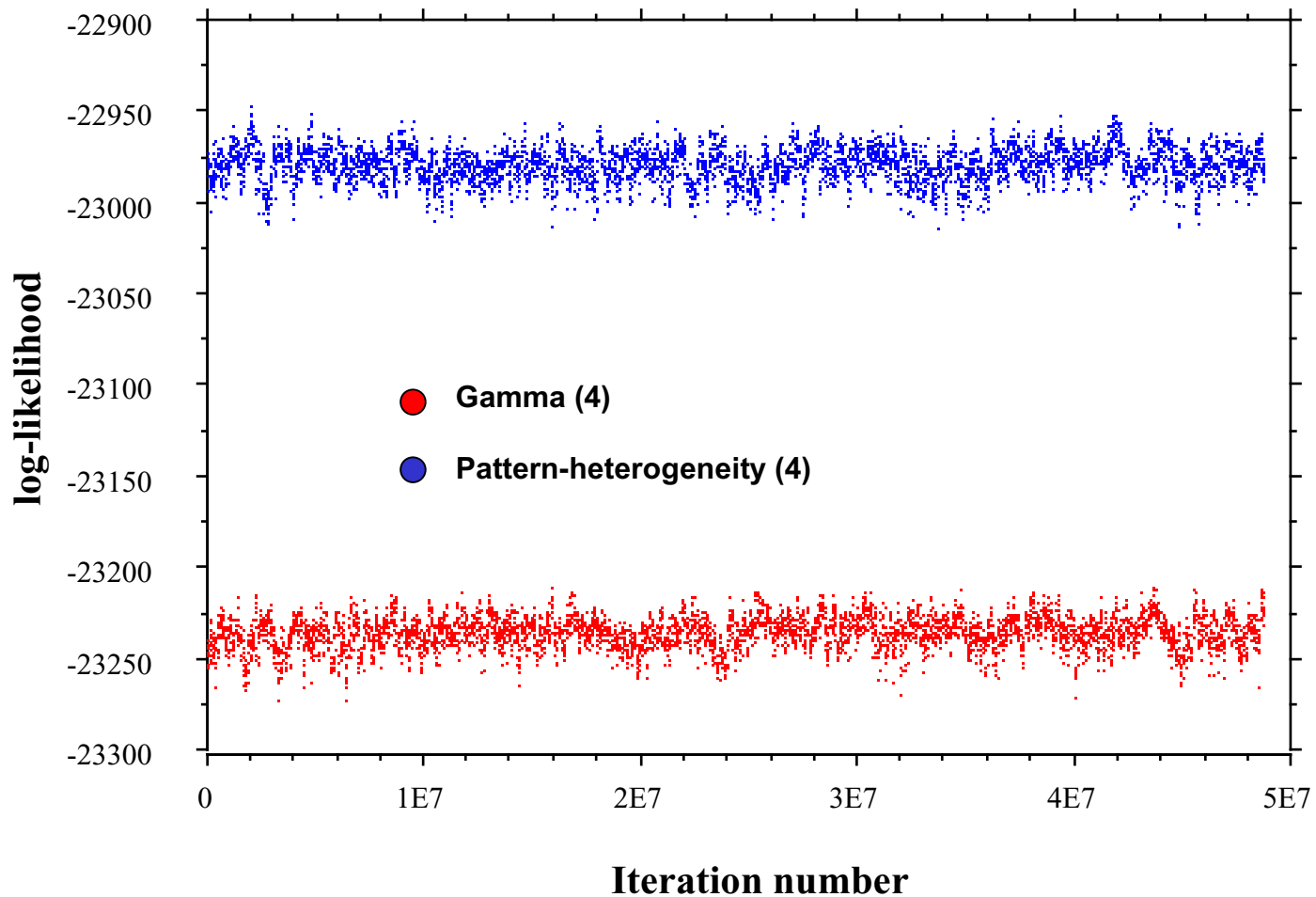
# Compensatory substitutions

# Detecting secondary structure: application of pattern-heterogeneity model to 12s mitochondrial DNA data on 57 mammals

**Detecting secondary structure; likelihoods of pattern-heterogeneity and Gamma models**

- ● Gamma (4)
- ● Pattern-heterogeneity (4)

log-likelihood

Iteration number

## Mammal 12S data: analysis of rate matrices

| | A <-> C | **A <-> G** | A <-> T | C <-> G | **C <-> T** | G <-> T | Tr/Tv |
|------|---------|---------|---------|---------|---------|---------|-------|
| Q1 | 13.75 | 10.52 | 7.89 | 1.69 | 59.65 | 3.34 | 5.26 |
| Q2 | 44.25 | 66.37 | 35.57 | 10.64 | 88.46 | 3.59 | 3.29 |
| **Q3** | **1.60** | **45.29** | **1.77** | **1.59** | **23.30** | **1.62** | **20.84** |
| Q4 | 0.30 | 0.80 | 0.18 | 0.20 | 1.62 | 0.10 | 6.26 |

**Best fit to loop or stem**

| | Q1 | Q2 | Q3 | Q4 |
|------|-----|-----|-----|-----|
| **Stem** | 57 | 13 | 122 | 272 |
| **Loop** | 102 | 150 | 57 | 236 |

Purines

Pyrimidines

**Conclusions**

Pattern-heterogeneity model can detect pattern-heterogeneity in simulated and real data and recover the parameters of the model of sequence evolution
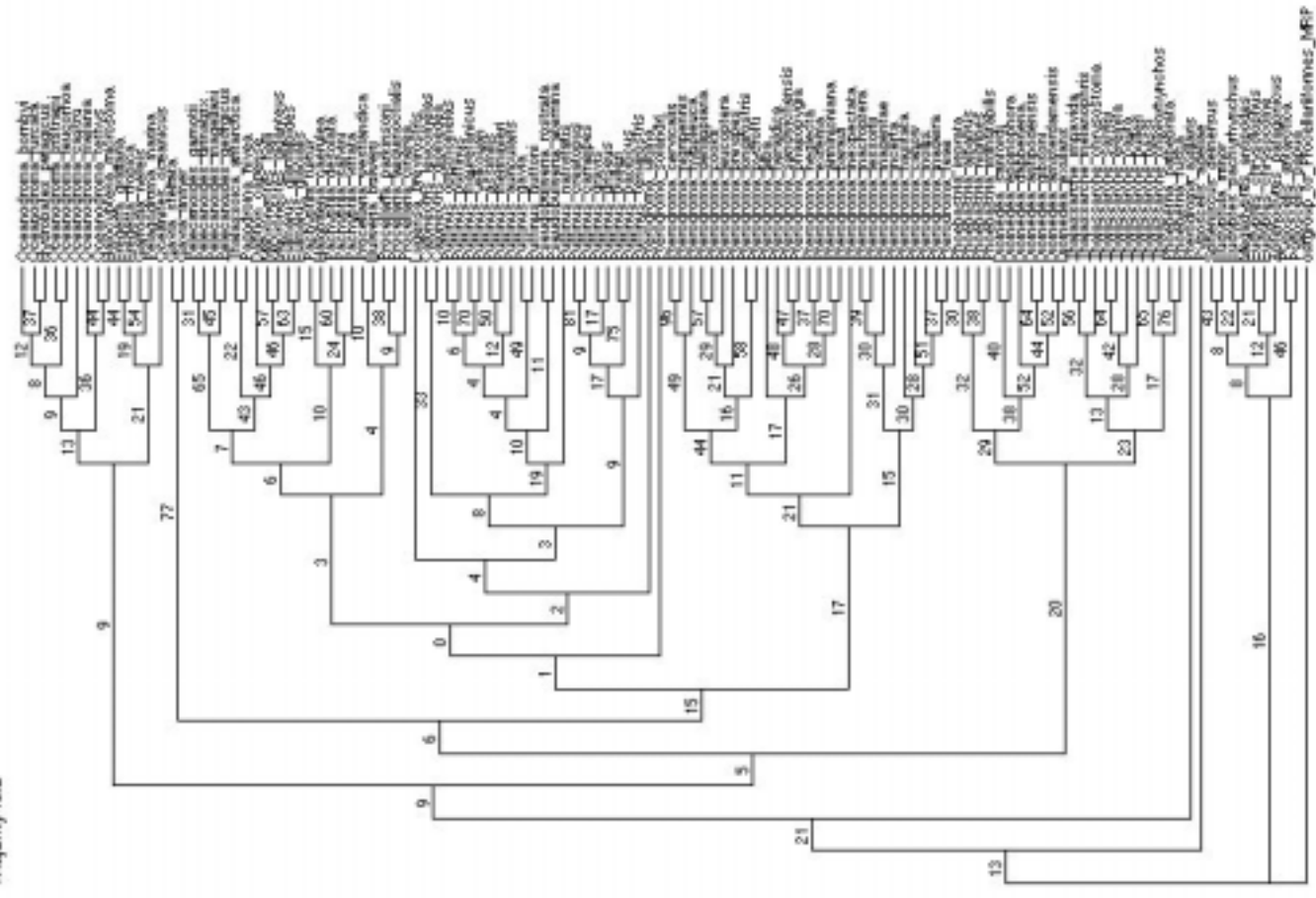
Can lead to large improvements in likelihood over homogeneous process model

Returns same likelihood as gamma-rates model for data with gamma rate variability and often improves upon gamma model in other situations

Can be used to investigate gene evolution (such as secondary structure) or be applied to concatenated data sets of multiple genes.

Software available from MP

Majority rule

Random tree of 60 tips used in all simulations