

# METABOLIC RANDOM FIELDS AND EVOLUTION

Stochastic models for the loss and gain of reactions in a metabolic network related by a phylogeny

# Contents

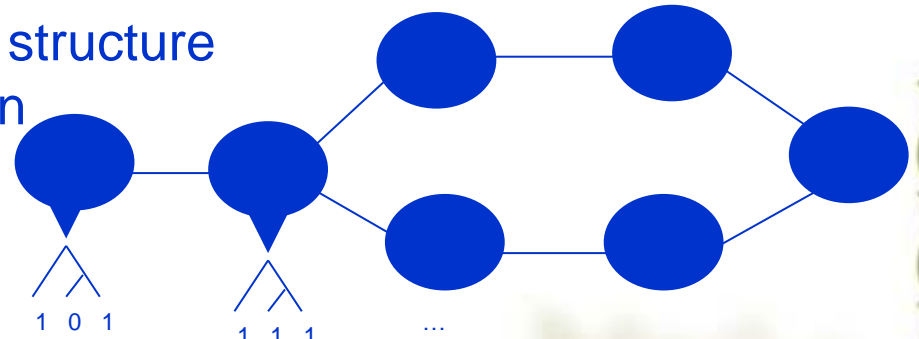
- ❖ Revision
- ❖ Sampling Hidden states,  $h$ 
  - MCMC
  - Wolff Algorithm
- ❖ Sampling Evolutionary Parameters,  $\lambda$  &  $\mu$
- ❖
- ❖ Data from KEGG: Extracting and Parsing
- ❖ Next Two Weeks
  - Data
  - General Method
- ❖ Limitations and Potential Extensions

# Revision

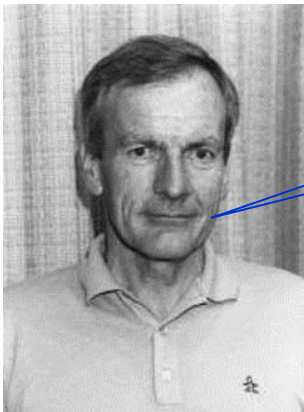
- ❖ Our **Aim**: to understand and model the gain and loss of reactions in metabolic networks for a known phylogeny
- ❖ Our **Approach**: Continuous time Markov Chain  
Felsenstein Algorithm  
Neighbour-dependent Model  
Markov Random Field  
  
and MORE...

# Sampling Hidden States, $h$

- Recall that we use a Markov Random Field (MRF) to model the dependency structure between hidden states at reaction sites



- We choose **Potts model** since it allows for any number of hidden states



$$p(h) = \exp \left\{ \beta \sum_{i \sim j} J(h_i, h_j) \right\} / Z(\beta)$$

where  $\beta$  is a measure of interaction strength between neighboring hidden states

and:

$$J(h_i, h_j) = \begin{cases} 1 & \text{if } h_i = h_j \\ 0 & \text{otherwise} \end{cases}$$

- By Bayes' rule, the posterior probability of hidden states  $h$  given the data  $\chi$  is given by:

$$\begin{aligned} p(h|\mathcal{X}) &= p(\mathcal{X}|h)p(h|\beta)/p(\mathcal{X}) \\ &= \sum_i p(\mathcal{X}|h_i)p(h|\beta)/Z(\mathcal{X}|\beta) \end{aligned}$$

- Under Potts model, we have:

$$p(h_i|h_{j \sim i}, \chi, \beta) \propto p(\chi_i|h_i)p(h_i|h_{j \sim i}) \propto p(\chi_i|h_i)e^{\beta \sum_{j \sim i} J(h_i, h_j)}$$

- Based on this, we can construct an MCMC scheme to sample  $h$ , with acceptance ratio:

$$r(h'_i|h_i) = \min \left\{ 1, \frac{p(h'_i|\chi)q(h_i|h'_i)}{p(h_i|\chi)q(h'_i|h_i)} \right\}$$

Unfortunately, such single-site proposal perform very badly in MRF!



WOLFF algorithm for the rescue

## ❖ **Wolff Algorithm** – *simple but efficient*

1. pick a reaction  $i$  at random, and form the sets  $C = \{i\}$ ,  $\bar{C} = \emptyset$
2. for every neighbor  $j$  of  $i$ , add  $j$  to  $C$  with probability  $1 - e^{-\beta J(h_i, h_j)}$
3. if  $j$  is added to  $C$ , all neighbors  $k$  of  $j$  are considered for inclusion in the cluster as in step 2, where  $k \notin C \cup \bar{C}$
4. after all reactions have been considered, set  $h_c = h'_i$  for all reactions  $c$  in  $C$  with probability  $\prod_{c \in C} \frac{p(\chi_c | h'_i)}{p(\chi_c | h_i)}$ , where  $h'_i$  is any state different from  $h_i$

In essence, if a cluster,  $C$  of reactions with hidden state  $h_i$  collectively has  $K$  neighbors in state  $h_i$  and  $K'$  neighbors in state  $h'_i$ , we've generated the cluster move from  $h_i$  to  $h'_i$  with probability  $(1 - e^{-\beta})^{|C|} e^{-\beta K}$  and from  $h'_i$  to  $h_i$  with probability  $(1 - e^{-\beta})^{|C|} e^{-\beta K'}$

As in the Potts model

$$\frac{p(h'_i|\chi)}{p(h_i|\chi)} = e^{\beta(K'-K)} \prod_{c \in C} \frac{p(\chi|h'_i)}{p(\chi|h_i)}$$

we see that the Wolff algorithm above generates exactly the same Metropolis-Hastings ratio as the single-site update, while the cluster update is much more likely to be accepted

# Sampling Rate Parameters, $\lambda$ & $\mu$

$\lambda$  and  $\mu$  are different for each hidden state  $h_k$ , denoted as  $\lambda_{h_k}$  and  $\mu_{h_k}$

But we assume these  $\lambda$  and  $\mu$  are all distributed *a priori* according to independent Gamma distributions,  $\text{gamma}(\alpha, \xi)$  (Mithani et al. 2010)

By Bayes' rule:

$$\begin{aligned} p(\lambda, \mu | \chi, h, t) &\propto p(\chi | h, \lambda, \mu) p(\lambda, \mu | h) \\ &= \prod_i p(\chi | h_i, \lambda_{h_i}, \mu_{h_i}) p(\lambda_{h_i} | \alpha, \xi) p(\mu_{h_i} | \alpha, \xi) \end{aligned}$$



Hence, we use MCMC to sample  $\lambda$  and  $\mu$ : a symmetric normal random walk on  $\log \lambda$  and  $\log \mu$  with Metropolis Hastings ratios:

$$r(\lambda'_{h^{(k)}} | \lambda_{h^{(k)}}) = \min\left\{1, \frac{f(\chi | \lambda'_{h^{(k)}}, \mu_{h^{(k)}})p(\lambda'_{h^{(k)}} | \alpha, \xi)}{f(\chi | \lambda_{h^{(k)}}, \mu_{h^{(k)}})p(\lambda_{h^{(k)}} | \alpha, \xi)}\right\}$$

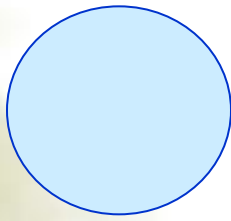
$$r(\mu'_{h^{(k)}} | \mu_{h^{(k)}}) = \min\left\{1, \frac{f(\chi | \mu'_{h^{(k)}}, \lambda_{h^{(k)}})p(\mu'_{h^{(k)}} | \alpha, \xi)}{f(\chi | \mu_{h^{(k)}}, \lambda_{h^{(k)}})p(\mu_{h^{(k)}} | \alpha, \xi)}\right\}$$

where  $f$  is the product of the likelihood of reactions in state  $h^{(k)}$  and can be calculated with the Felsenstein algorithm

We have looped these sampling algorithms together and tested it on simulated datasets.

# “SPAM metabolism network”

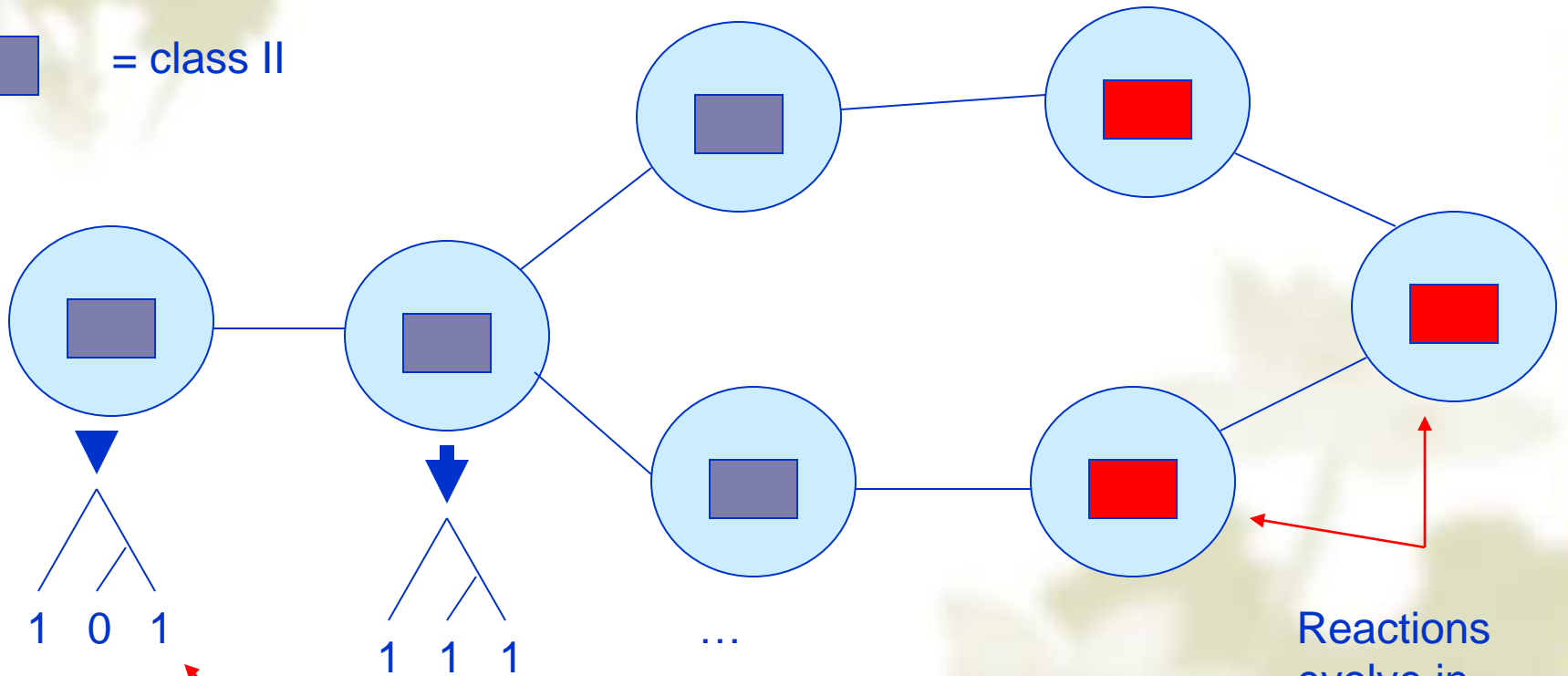
reaction =



= class I



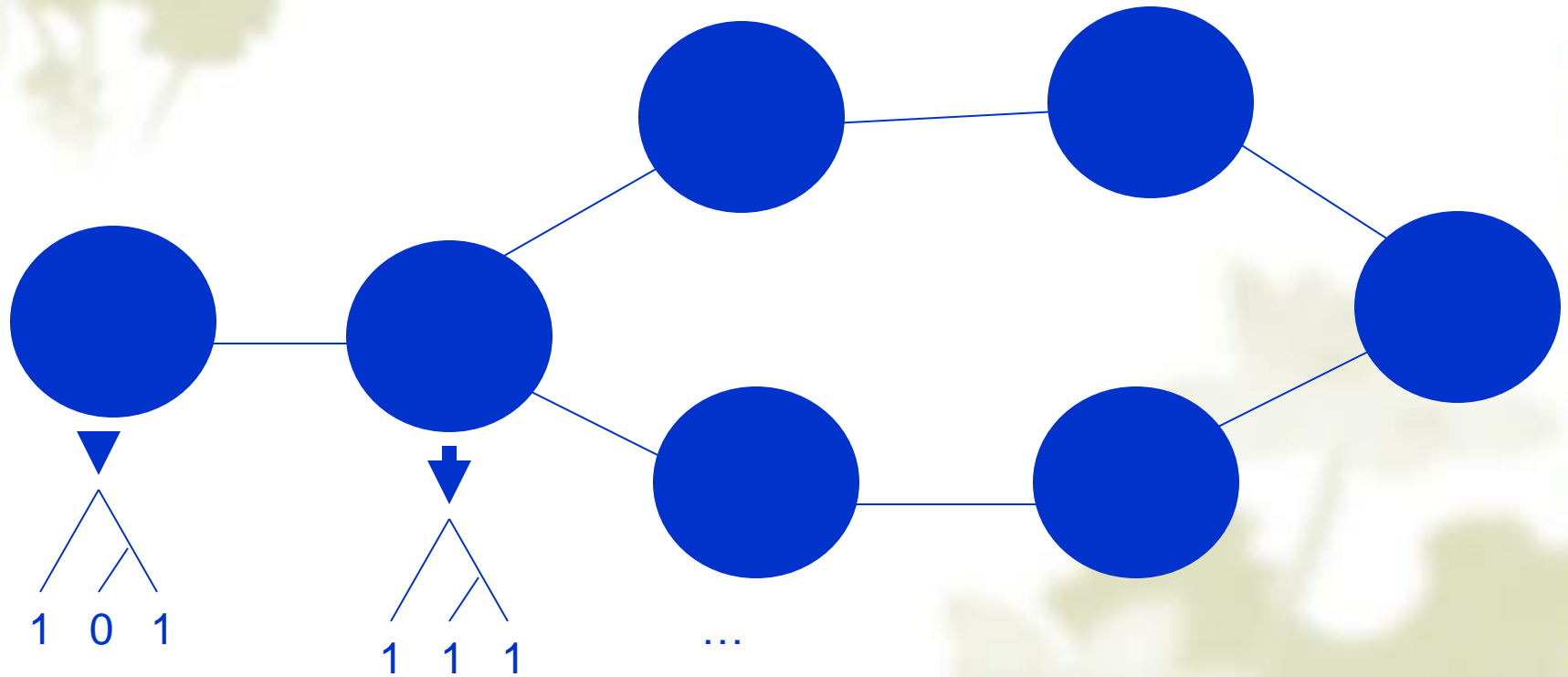
= class II



Reactions evolving at rate  $Q$

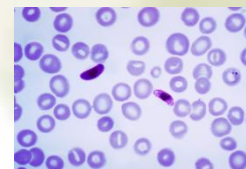
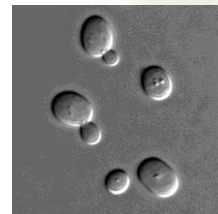
Reactions evolve in accord with  $Q$

# SPAM metabolism network



Network hidden: cannot see underlying state, but we know the phylogenies

# Data



# Data

- ❖ rn: 147
- ❖ hsa: 87
- ❖ dme: 86
- ❖ ath: 95
- ❖ sce: 72
- ❖ pfa: 57
- ❖ eco: 89

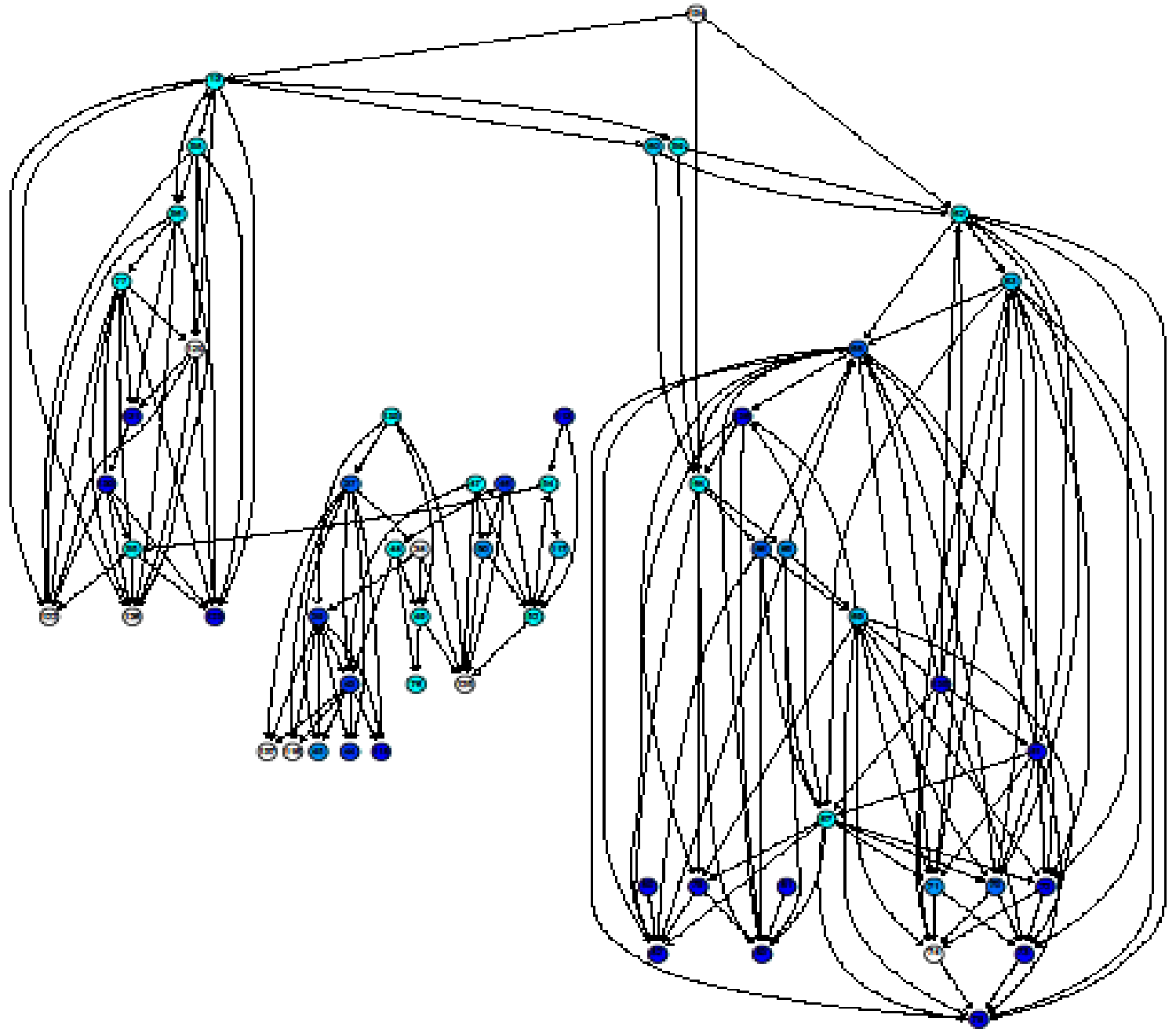
- ❖ rn00010 (glycolysis)
- ❖ rn00020 (TCA cycle)
- ❖ rn00030 (pentose phosphate)
- ❖ rn00040 (pentose/glucuronate interconversions)
- ❖ rn00051 (fructose/mannose)
- ❖ rn00052 (galactose)
- ❖ rn00053 (ascorbate/aldarate)
- ❖ rn00061 (fatty acid biosynthesis)

# Data

- ❖ Load maps into R using *KEGGgraph*
- ❖ Parse reactions
- ❖ Save as edgelists and nodelists (load into simulator later)



# Data

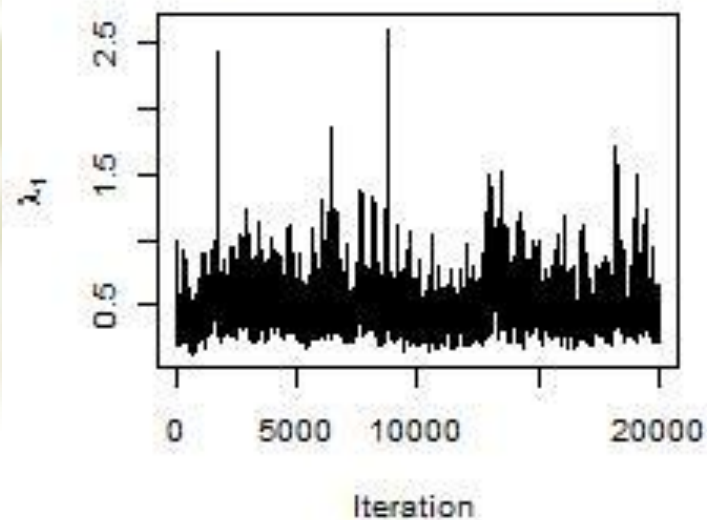


# Results

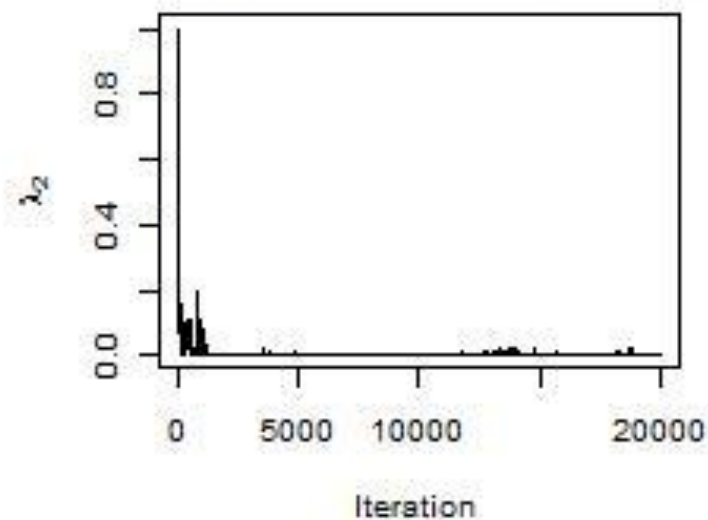
- ❖ For now 1 network with 6 species phylogeny.
- ❖ Branch lengths of phylogenetic tree are randomly generated.
- ❖ Each reaction takes two possible annotations or states. These correspond to the rates of addition or deletion of reactions in the network.
- ❖ Run MCMC on  $\lambda$   $\mu$  and  $h$  (the hidden state).



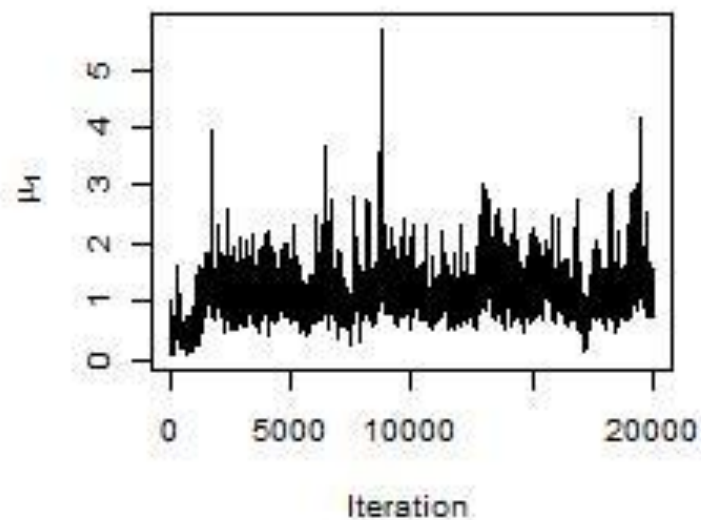
MC\_aa1.1\_b0.1\_K2\_cl1\_lambda\_1



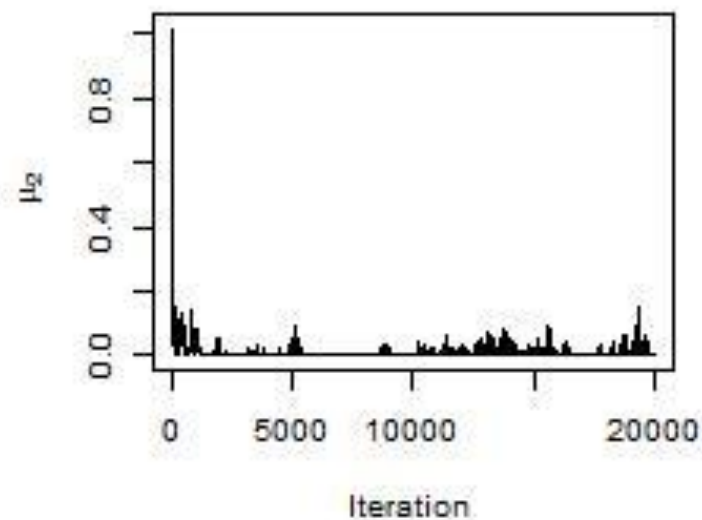
MC\_aa1.1\_b0.1\_K2\_cl1\_lambda\_2



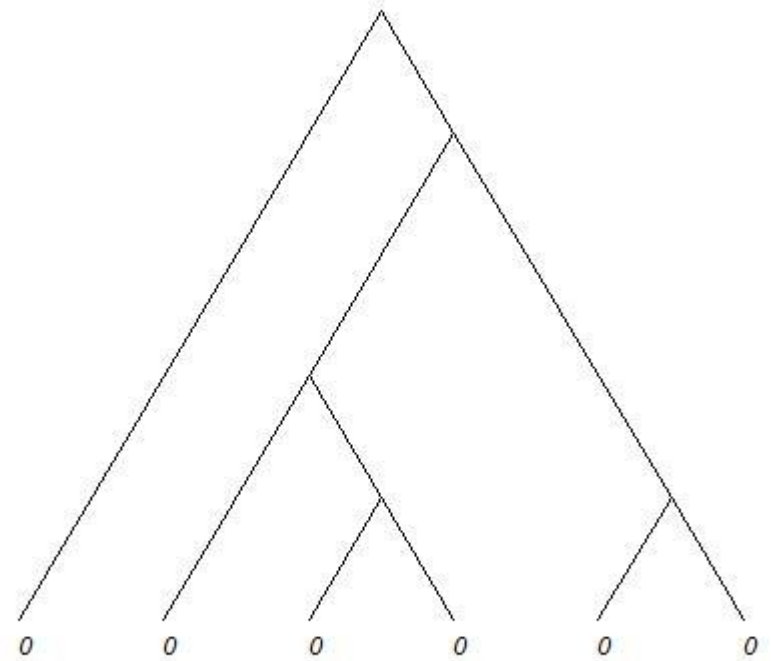
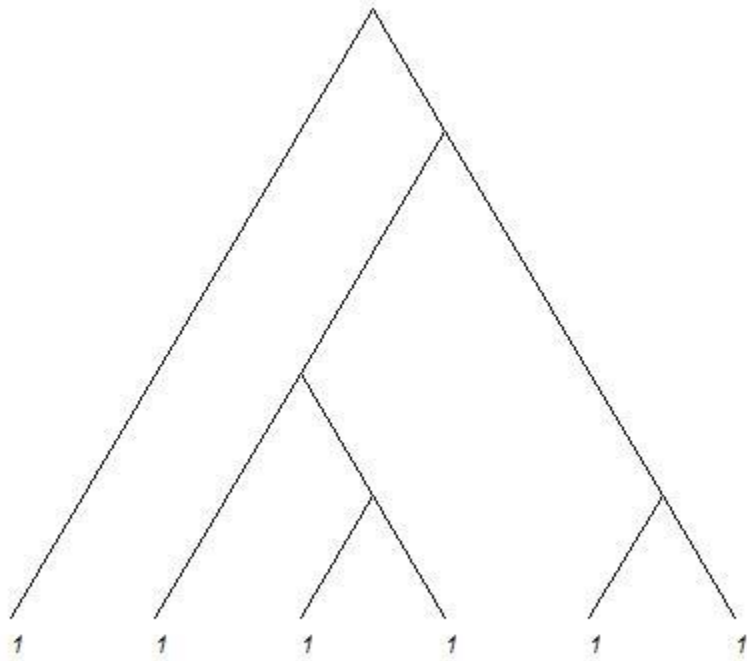
MC\_aa1.1\_b0.1\_K2\_cl1\_mu\_1



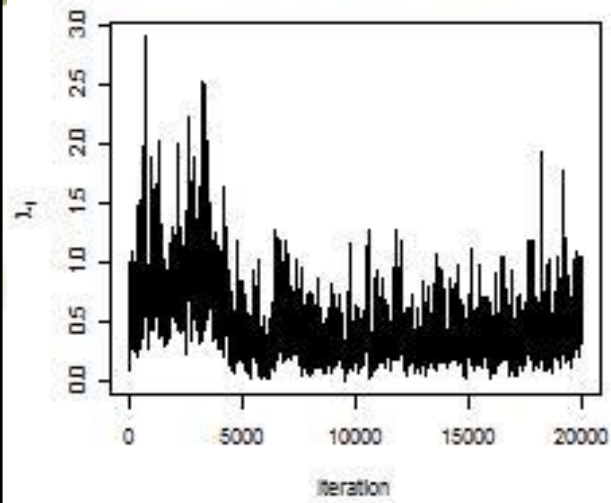
MC\_aa1.1\_b0.1\_K2\_cl1\_mu\_2



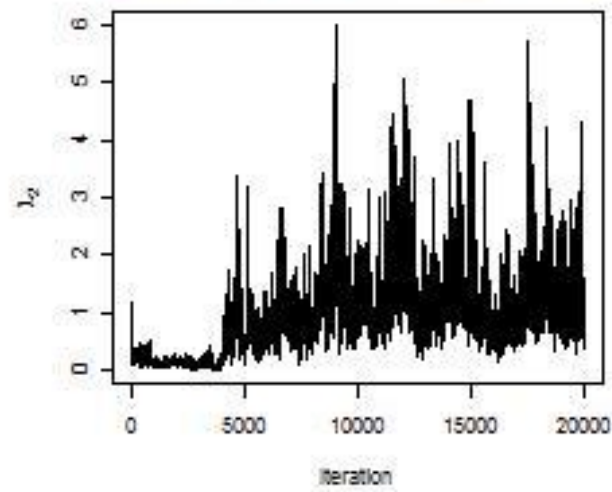
# Conserved Reactions



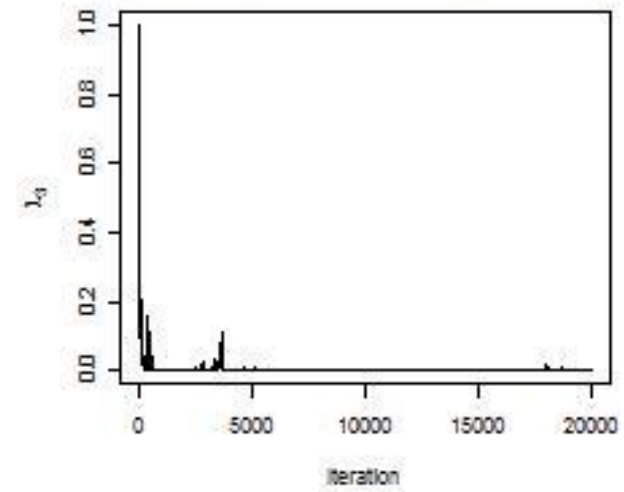
MC\_aa1.1\_b0.2\_K3\_cl1\_lambda\_1



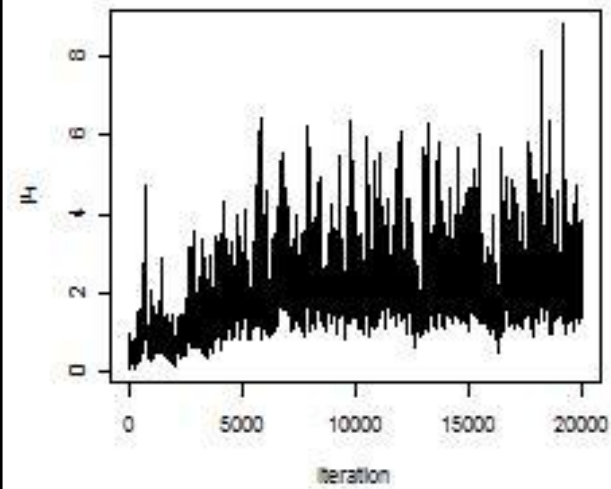
MC\_aa1.1\_b0.2\_K3\_cl1\_lambda\_2



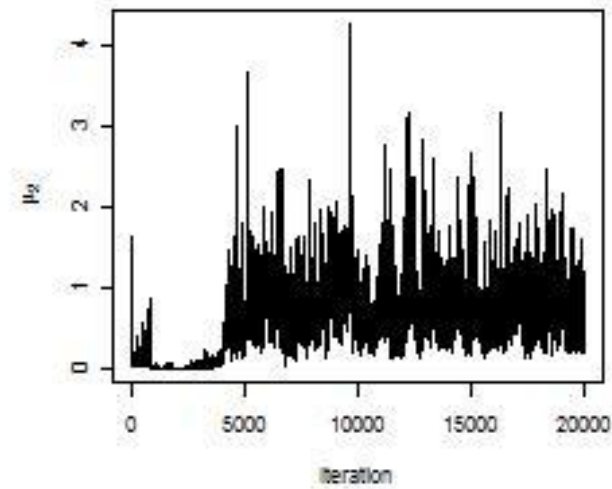
MC\_aa1.1\_b0.2\_K3\_cl1\_lambda\_3



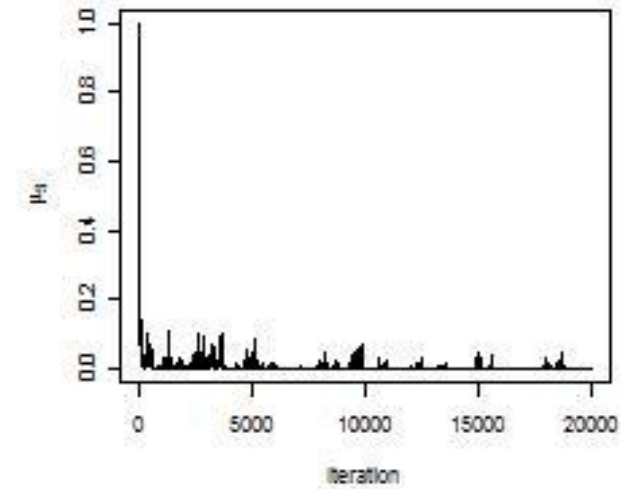
MC\_aa1.1\_b0.2\_K3\_cl1\_mu\_1



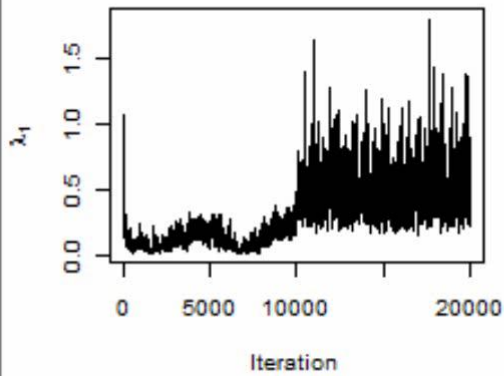
MC\_aa1.1\_b0.2\_K3\_cl1\_mu\_2



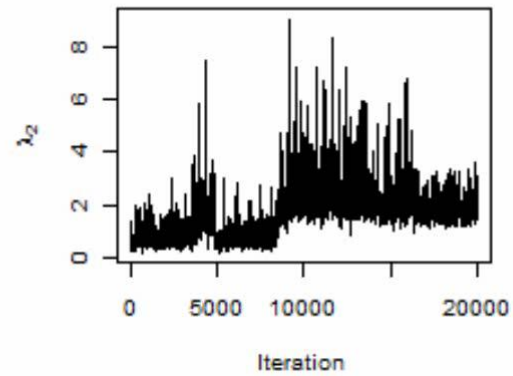
MC\_aa1.1\_b0.2\_K3\_cl1\_mu\_3



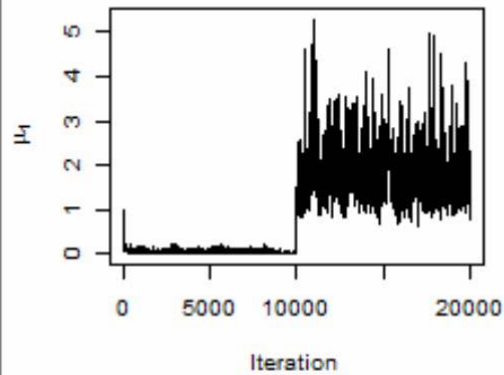
MC\_aa1.1\_b0\_K2\_ci0\_lambda\_1



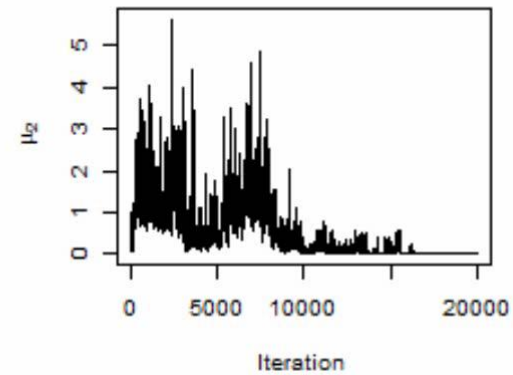
MC\_aa1.1\_b0\_K2\_ci0\_lambda\_2



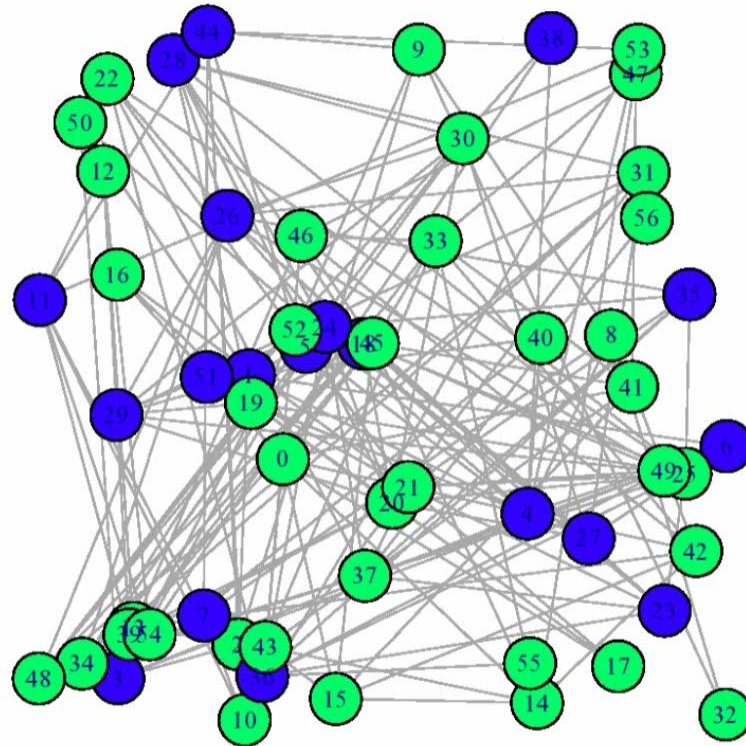
MC\_aa1.1\_b0\_K2\_ci0\_mu\_1



MC\_aa1.1\_b0\_K2\_ci0\_mu\_2



iteration: 1000 , K = 2



# Why Such Variability??

- ❖ Variation of interaction strength between neighbouring reactions.
- ❖ Variation in jump sizes.
- ❖ Cluster flips or single state flips.
- ❖ Number or possible hidden states.
- ❖ Prior for addition / deletion probabilities.
- ❖ Parameterisation of these?

# What's Next

- ❖ Account for popularity tendency.
- ❖ Testing for convergence.
- ❖ Different species? more species?
- ❖ Branch lengths?
- ❖ Enough nodes? (could combine networks)
- ❖ Write the report



**THE  
END**