

# Phylogenetic Graphical Models and RevBayes: Introduction

Fred(rik) Ronquist  
Swedish Museum of Natural History,  
Stockholm, Sweden

# Statistical Phylogenetics

- Statistical approaches increasingly important:
  - Difficult problems requiring accurate and unbiased inference (e.g., structure of rapid radiations)
  - More aspects of molecular evolution being examined (structural dependencies, etc)
  - Combination of background knowledge and sequence information (e.g., divergence time estimation)
- Modeling explosion, especially in the Bayesian context
- Challenging for empiricists to communicate and correctly understand models
- Challenging for developers of inference software



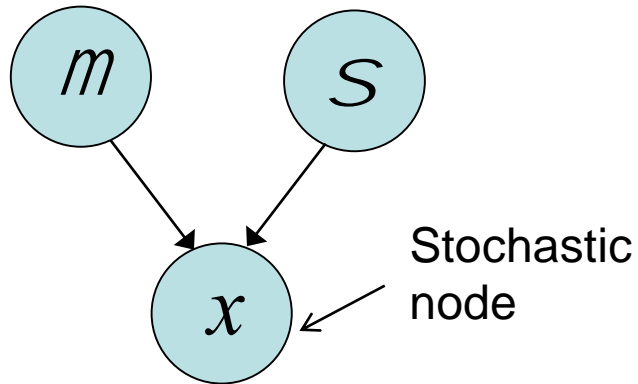
# Probabilistic Graphical Models

- Theoretical framework for specifying dependencies in complex statistical models
- Allows a complex model to be broken down into conditionally independent distributions
- Closely related to standard statistical model formulae:

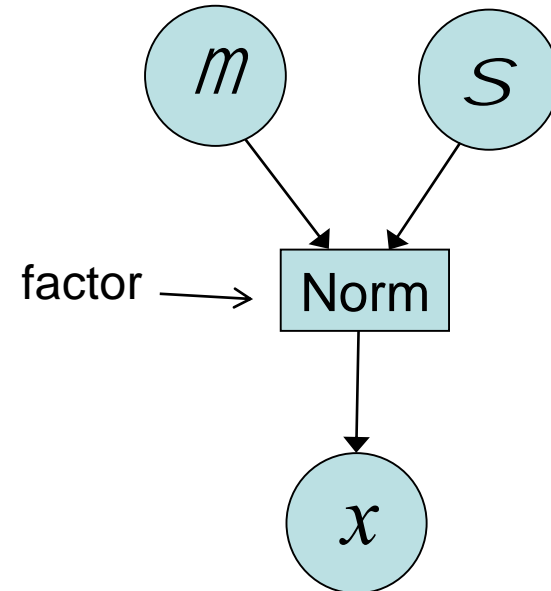
$$x \sim \text{Norm}(\mu, \sigma)$$

- Extensive literature on generic algorithms that apply to model graphs

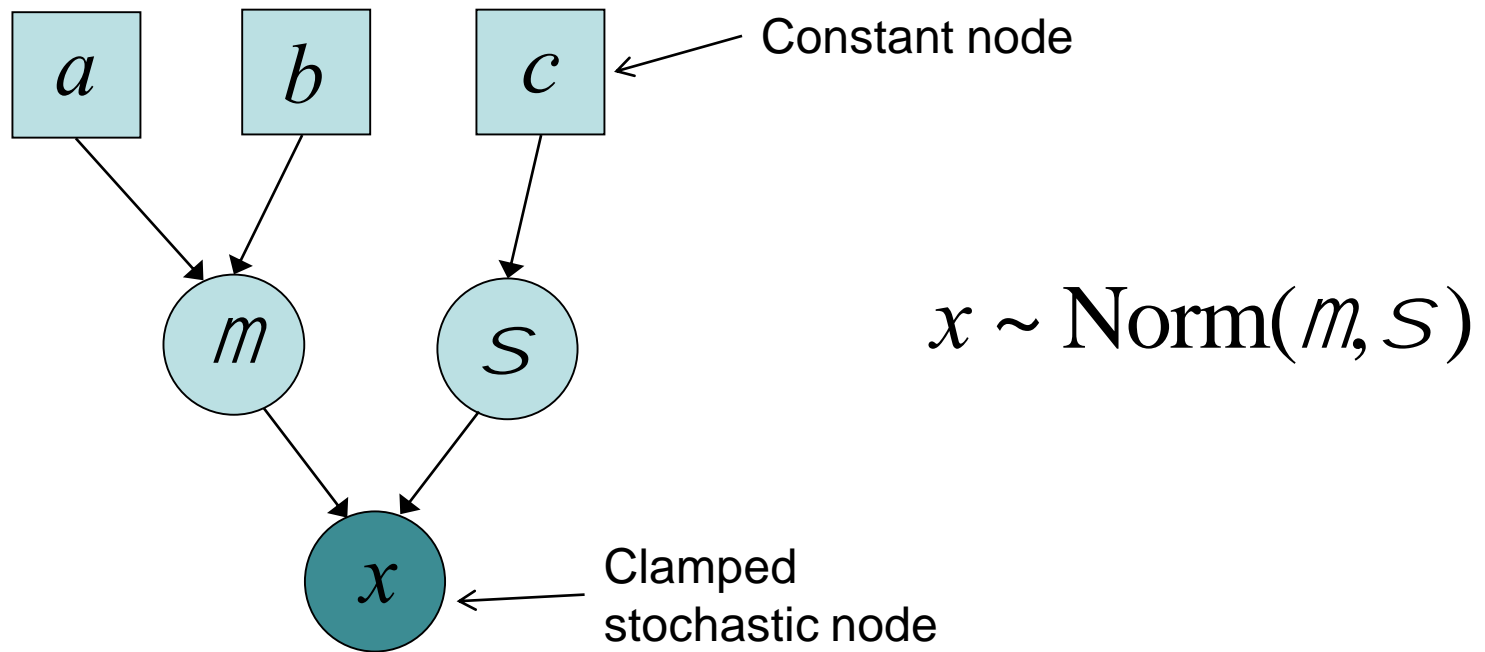
$$x \sim \text{Norm}(m, S)$$



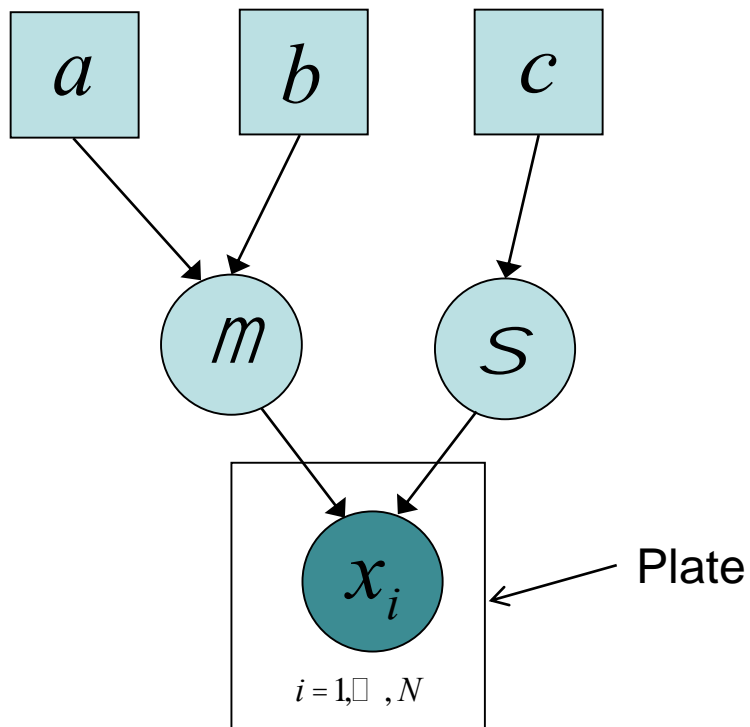
Graphical Model  
Compact Form



Factor Graph



Hierarchical Graphical Model



$$x \sim \text{Norm}(m, S)$$

Hierarchical Graphical Model

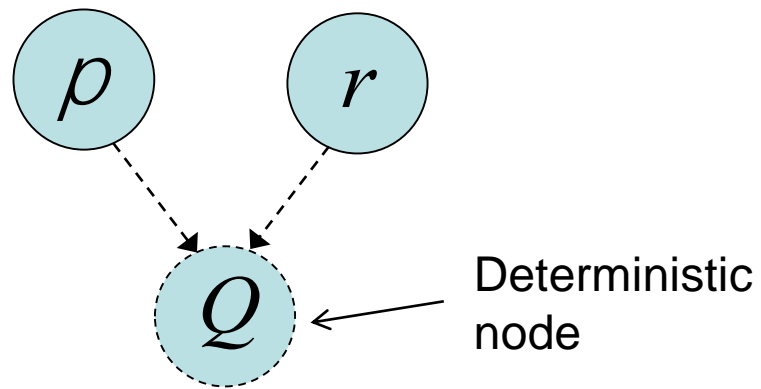
$$Q = \begin{pmatrix} - & \rho_C r_{AC} & \rho_G r_{AG} & \rho_T r_{AT} \\ \rho_A r_{AC} & - & \rho_G r_{CG} & \rho_T r_{CT} \\ \rho_A r_{AG} & \rho_C r_{CG} & - & \rho_T r_{GT} \\ \rho_A r_{AT} & \rho_C r_{CT} & \rho_G r_{GT} & - \end{pmatrix}$$

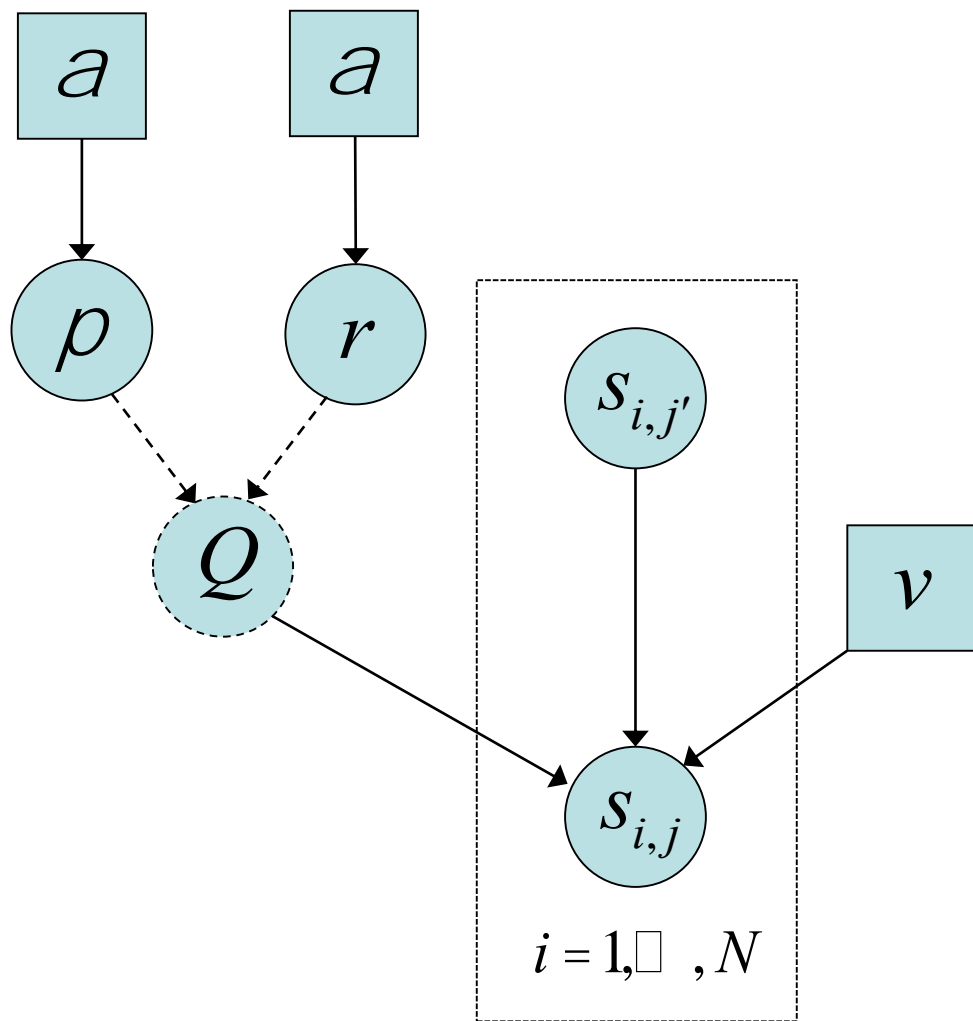
General Time Reversible  
(GTR) substitution model

$\rho$  Stationary state frequencies

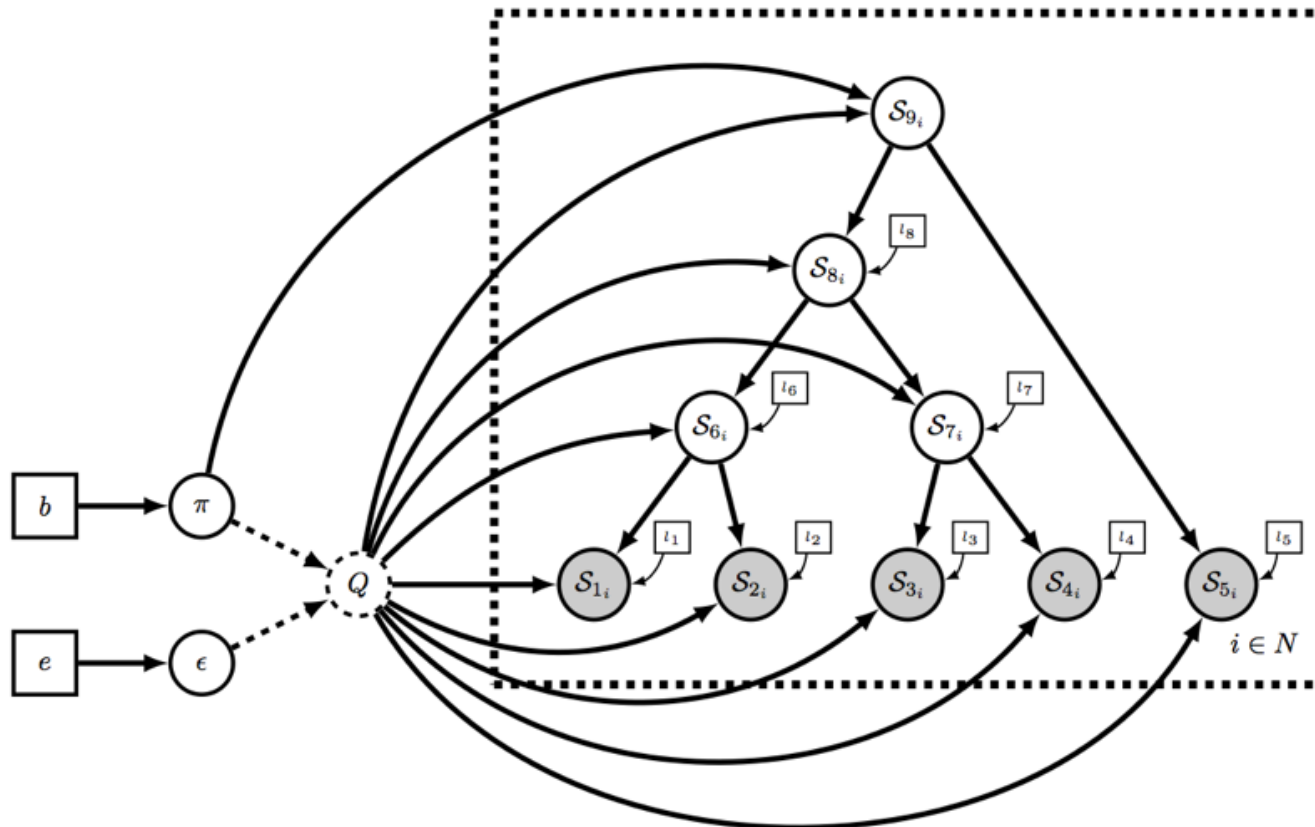
$r$  Exchangeability rates



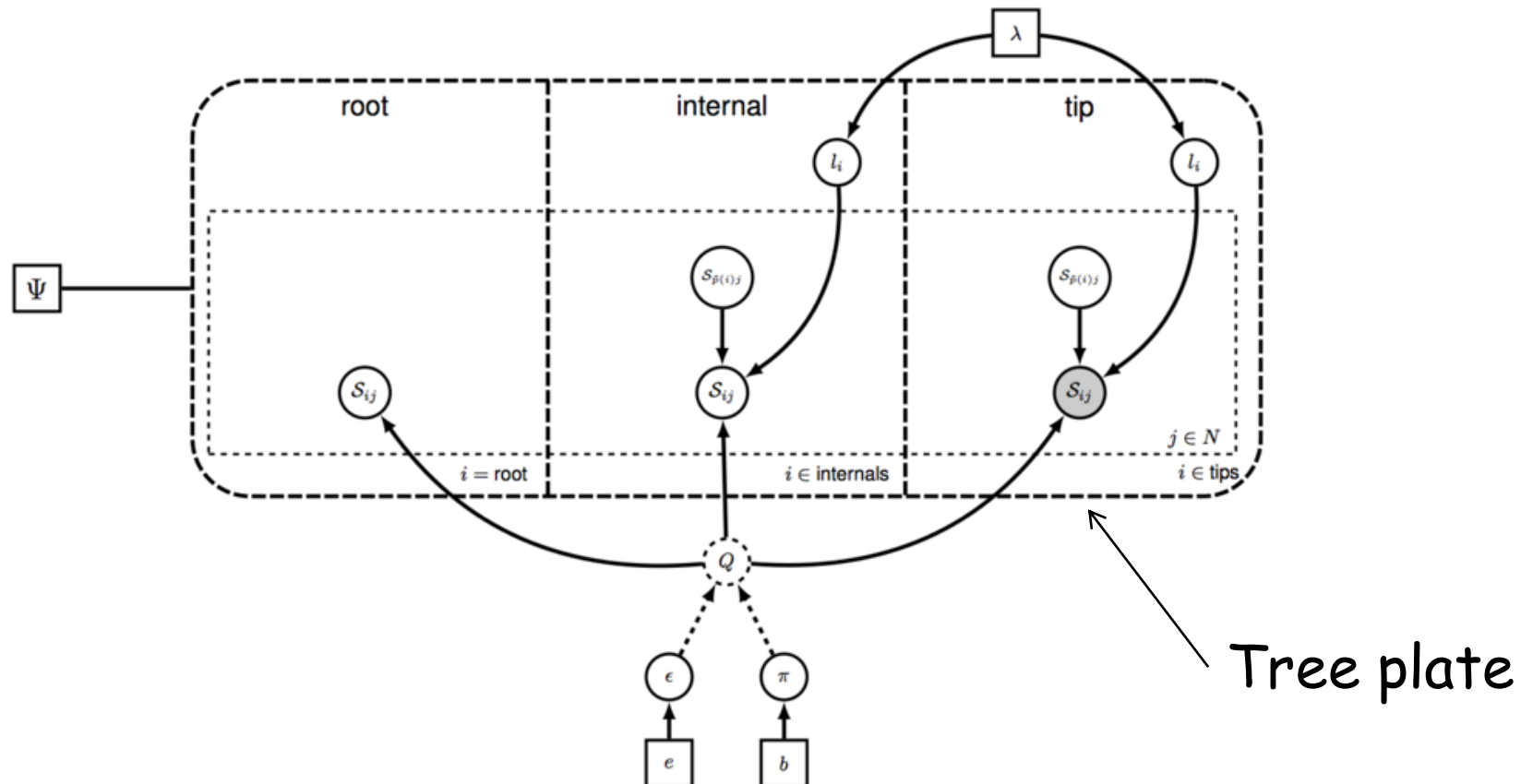


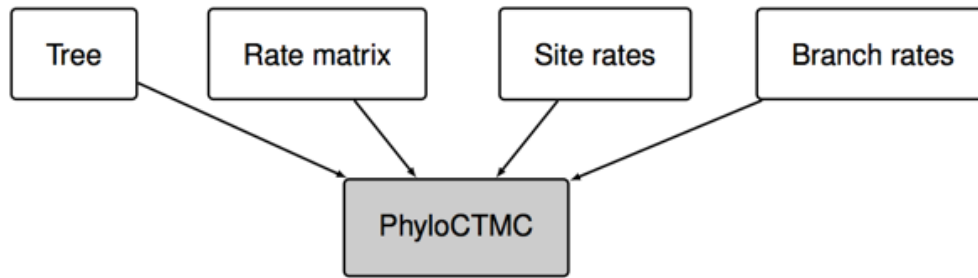


# GTR Phylogeny Model

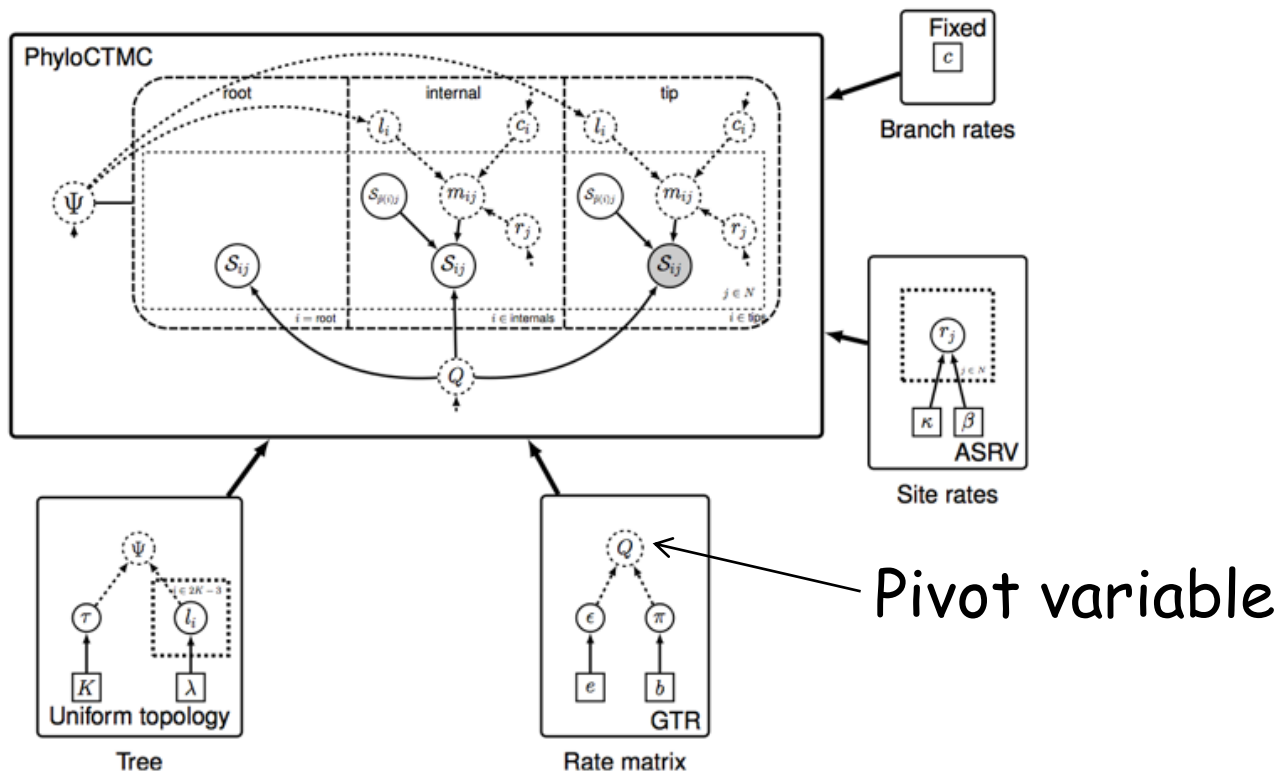


# Tree Plate Representation





# Modular Representation



# RevBayes Project

- Interactive computing environment intended primarily for Bayesian phylogenetic inference
- Uses a special language, Rev, for constructing probabilistic phylogenetic and evolutionary graphical models interactively, step by step
- Rev is similar to R and the BUGS modeling language
- RevBayes provides generic computing machinery for simulation, inference and model testing



# Basic properties of the Rev language

# There are three kinds of statements in the language

# 1. Arrow assignment (value assignment, create constant nodes)

```
> a <- 4          # Give a the value 4
> b <- sqrt(a)     # Give b the value of sqrt(a), that is, 2
> b               # Print the value of b
2
```

# 2. Equation assignment (create deterministic nodes)

```
> c := sqrt(a)     # Make c a dynamic function node evaluating sqrt(a)
> c
2
> a <- 9           # Give a the value 9
> b               # Print the value of b
2
> c               # Print the value of c
3
```



# Basic properties of the Rev language

# 3. Tilde assignment (create stochastic variables (nodes))

> a ~ dnExp( rate = x )      # a is drawn from exp dist with rate = x

# Basic properties of the Rev language

```
# -----
```

```
# Declaring and defining functions
```

```
# -----
```

```
> function foo ( x ) { x * x }
```

```
> foo( 2 )
```

```
4
```

```
# If you wish, you can specify types as well
```

```
> function PosReal foo ( Real x ) { x * x }
```

```
# Without explicit types, RevObject is the assumed type
```

```
# -----
```

```
# Declaring and defining new types
```

```
# -----
```

```
> class myclass : Move {
```

```
+   Real myTuningParam;
```

```
+   procedure Real move( Real x ) { myTuningParam * x }
```

```
+ }
```

```
# Inheritance, function overriding and overloading
```

# A complete MCMC analysis in Rev

```
a <- -1.0
```

```
b <- 1.0
```

```
mu ~ dnUnif(a, b)
```

```
sigma ~ dnExp(1.0)
```

```
for (i in 1:10) {
```

```
  x[i] ~ dnNorm(mu, sigma)
```

```
  x[i].clamp(0.5)
```

```
}
```

```
mymodel = model(mu) # Any stochastic node in the model works
```

```
mymcmc = mcmc(mymodel)
```

```
mymcmc.run(1000)
```

```

# definition of the myGTR function ("Ziheng's favorite")
function model myGTR (CharacterMatrix data) {

  # describe Q matrix
  pi ~ dflatdir(4);
  r ~ dflatdir(6);
  Q := gtr(pi, r);

  # describe tree
  tau ~ dtopuni(data.taxa(), rooted=false);

  # gamma shape
  alpha ~ dunif(0.0, 50.0);

  # discrete gamma mixture
  for (i in 1:4)
    catRate[i] := qgamma(i*0.25-0.125, alpha, alpha);
  for (i in 1:data.size())
    ratecat[i] ~ dcat(simplex(0.25,0.25,0.25,0.25));

  # associate distributions with tree parts
  for (i in 1:data.size()) {
    for (n in 1:tau.numNodes()) {
      if (tau.isTerminal(n)) {
        tau.length[n] ~ exp(1.0);
        tau.state[n] ~ ctmc(Q, e.length*catRate[ratecat[i]],
          tau.state[tau.parent(n)]);
        tau.state[n] <- data[i][tau.tipIndex(n)];
      }
      else {
        tau.length[n] ~ exp(10.0);
        tau.state[n] ~ ctmc(Q, e.length*catRate[ratecat[i]],
          tau.state[n]);
      }
    }
  }

  # return model
  return model( Q );
}

```

Definition of  
a new  
phylogenetic  
model

Appr. 20 lines

# Complexity hidden from normal user

```
# Read in data
myData <- read( "data.nex" )

# Apply model
myModel = zihengGTR( myData )

# Construct mcmc
myMCMC = mcmc( myModel )

# Run mcmc
myMCMC.run(10000)
```

# RevBayes

- RevBayes is still experimental software
- Help is still incomplete
- There may be various bugs and other problems, for instance related to type conversion
- A number of practical features are still missing
  - Post-hoc analysis very limited
  - No support for convergence diagnostics on the fly

# RevBayes Limitations

- Coarse-grained representation of tree plates -> some models are impossible to specify
- No ability to manipulate modules
- No default moves or monitors
- Some programming features are missing, notably specification of user-defined types
- Most additions need to be made in the back end using heavily templated C++ code

# RevLang Project

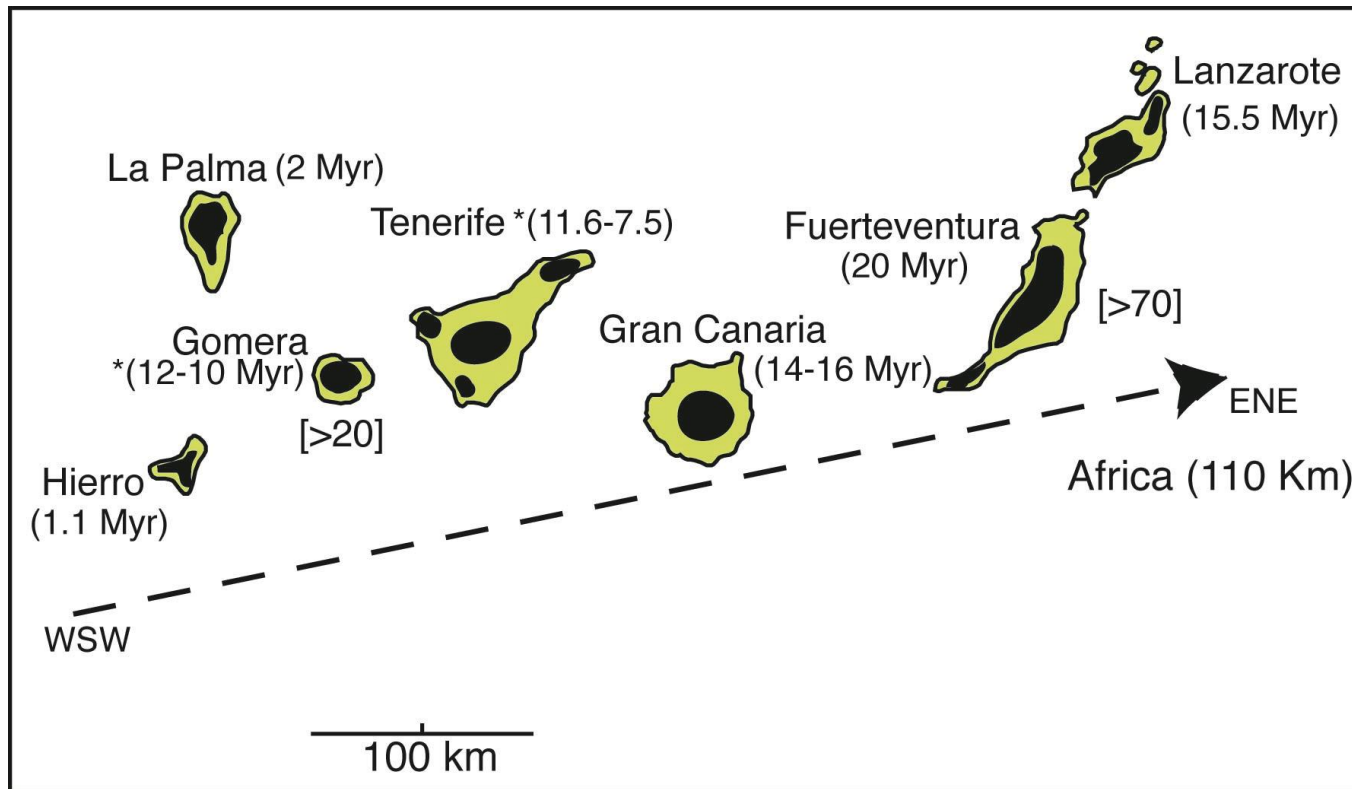
- Rev -> probabilistic programming language
- Moves, distributions etc programmed in Rev itself
- Professional JIT compiler
- Fully supported interactive environment, like R
- Modular design makes it easy to extend the environment



# RevBayes Exercises

- Download and install RevBayes according to instructions at the RevBayes web site (<http://revbayes.com>).
- Download material and follow tutorials of interest
- RevBayes code at github:  
<https://github.com/revbayes/revbayes>

# The Canary Islands



# CANARY ISLANDS (Islas Canairas)

(Spain)

1:6,000,000

*Ilhas Selvagens*  
(Portugal)

30°

*Alagranza*

*Isla Graciosa*

*Lanzarote*

Arrecife

Puerto de Rosario

*Fuerteventura*

Las Palmas  
de Gran Canaria

*Gran Canaria*

MOROCCO

LAAYOUNE

WESTERN SAHARA

14°

16°

18°

28°

*La Palma*

Los Llanos  
de Aridane

Santa Cruz  
de la Palma

Santa Cruz  
de Tenerife

*Tenerife*

Pico de Teide  
3718

*Isla de la Gomera*

San Sebastián  
de la Gomera

Frontera

*Hierro*

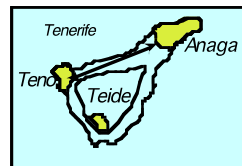
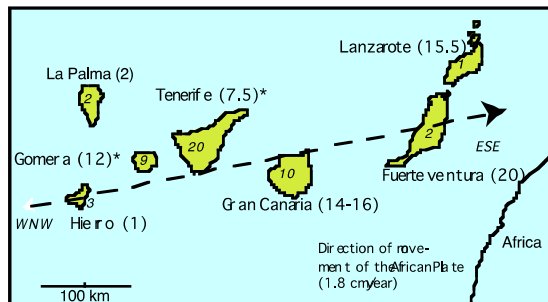




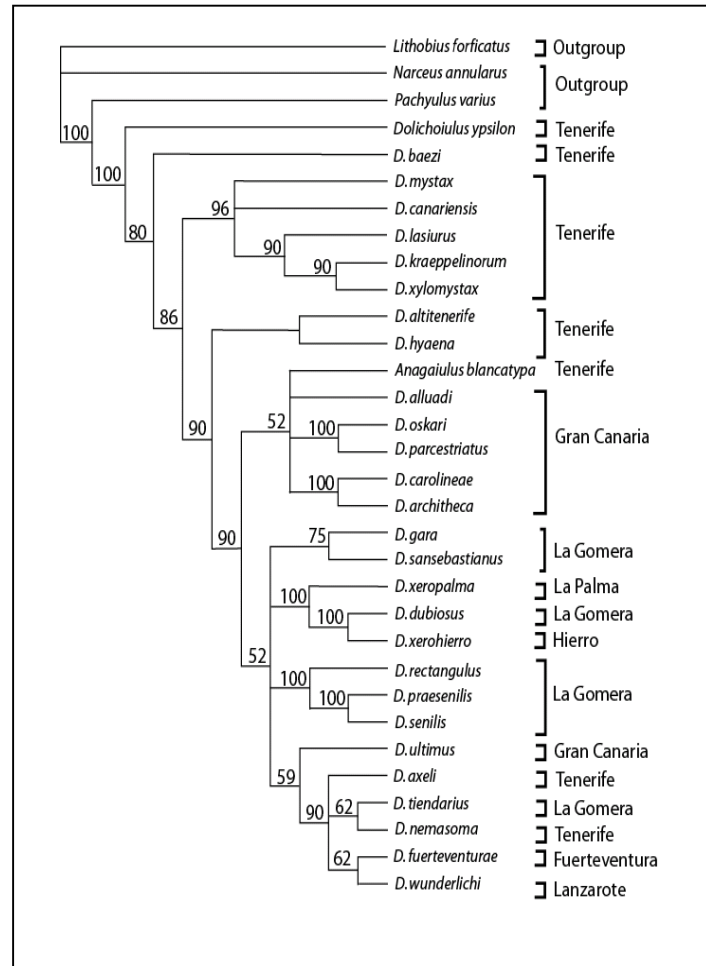
# Dolichoziulus (Diplopoda)



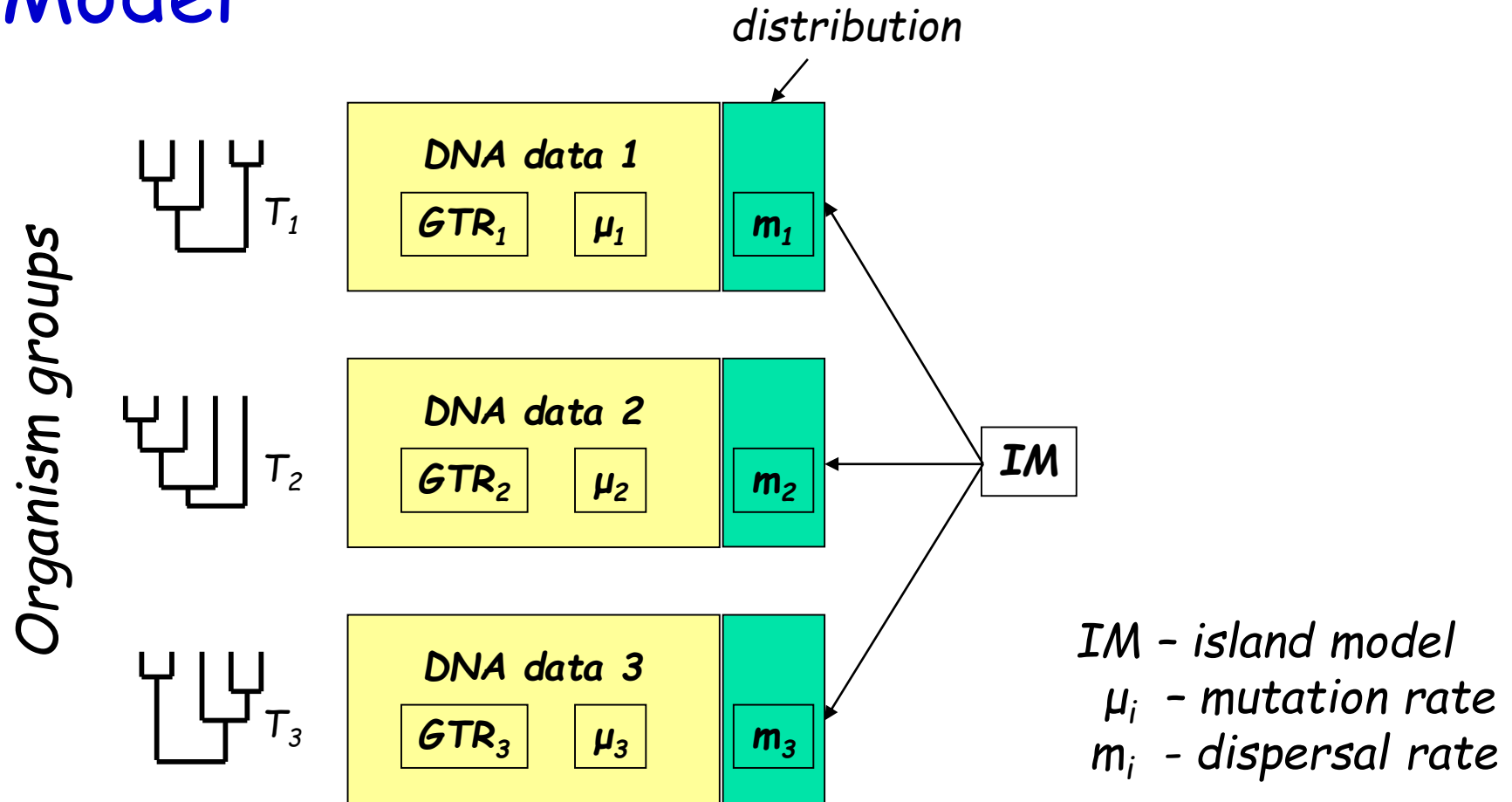
*Dolichoziulus* (Diplopoda, Julida, Julidae, Pachyulinae)



46 endemic species



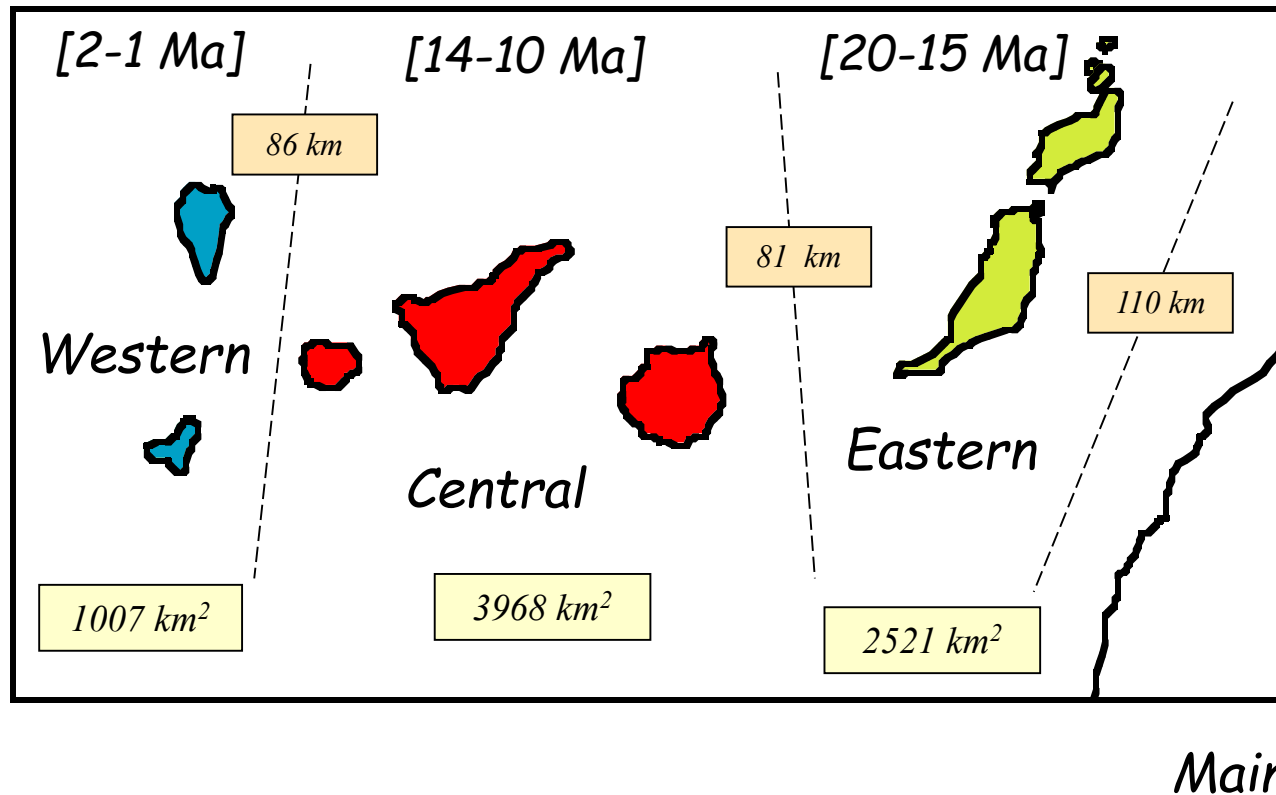
# Model



# Inference

Bayesian inference using MCMC sampling,  
accommodating uncertainty in all model parameters

# Canary Islands: 3-island model



# Instantaneous rate matrix

to

from

$$Q = \begin{pmatrix} & [A] & [B] & [C] & [D] \\ [A] & - & \pi_B r_{AB} & \pi_C r_{AC} & \pi_D r_{AD} \\ [B] & \pi_A r_{AB} & - & \pi_C r_{BC} & \pi_D r_{BD} \\ [C] & \pi_A r_{AC} & \pi_B r_{BC} & - & \pi_D r_{CD} \\ [D] & \pi_A r_{AD} & \pi_B r_{BD} & \pi_C r_{CD} & - \end{pmatrix}$$

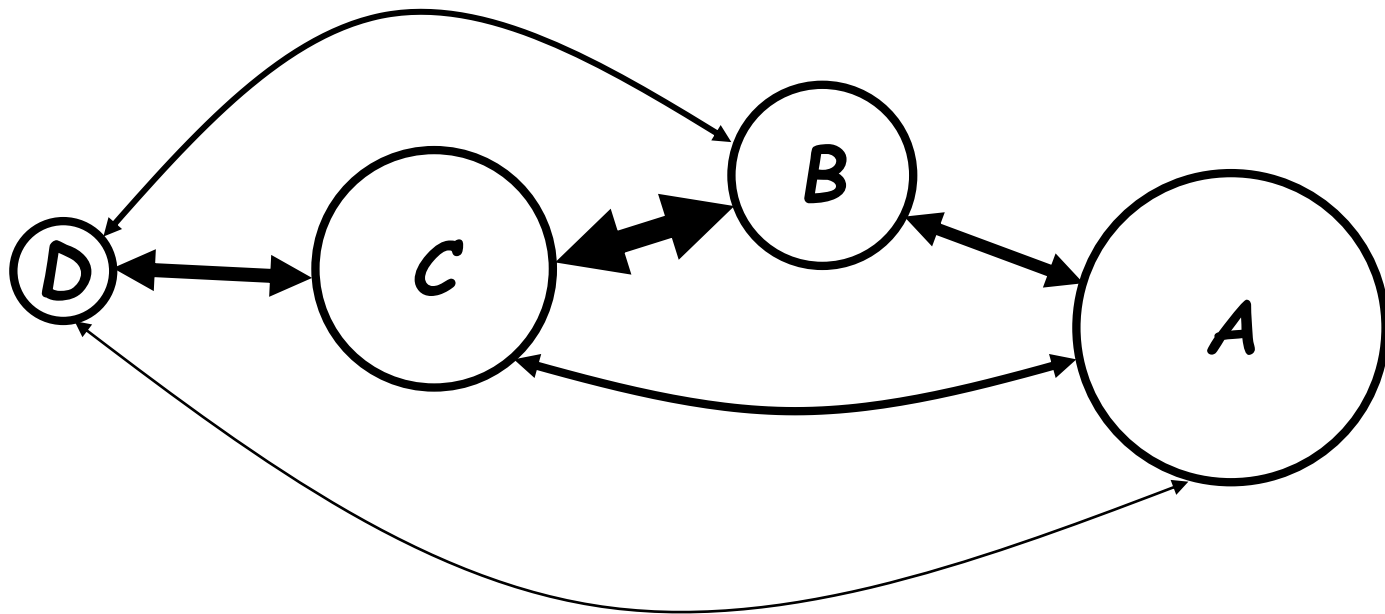
$\pi_i$  Relative carrying capacity of island  $i$

$r_{ij}$  Relative dispersal rate between islands  $i$  and  $j$



# Island "GTR" model

*Time reversible continuous time Markov chain*



*Differing relative carrying capacities of islands*  
*Differing intensities of biotic exchange between islands*



*Laura Martinez*

## Canary Island flora

*15 groups, 567 lineages*



*Javier Fuertes*







*Raquel Martin*

# Canary Island fauna

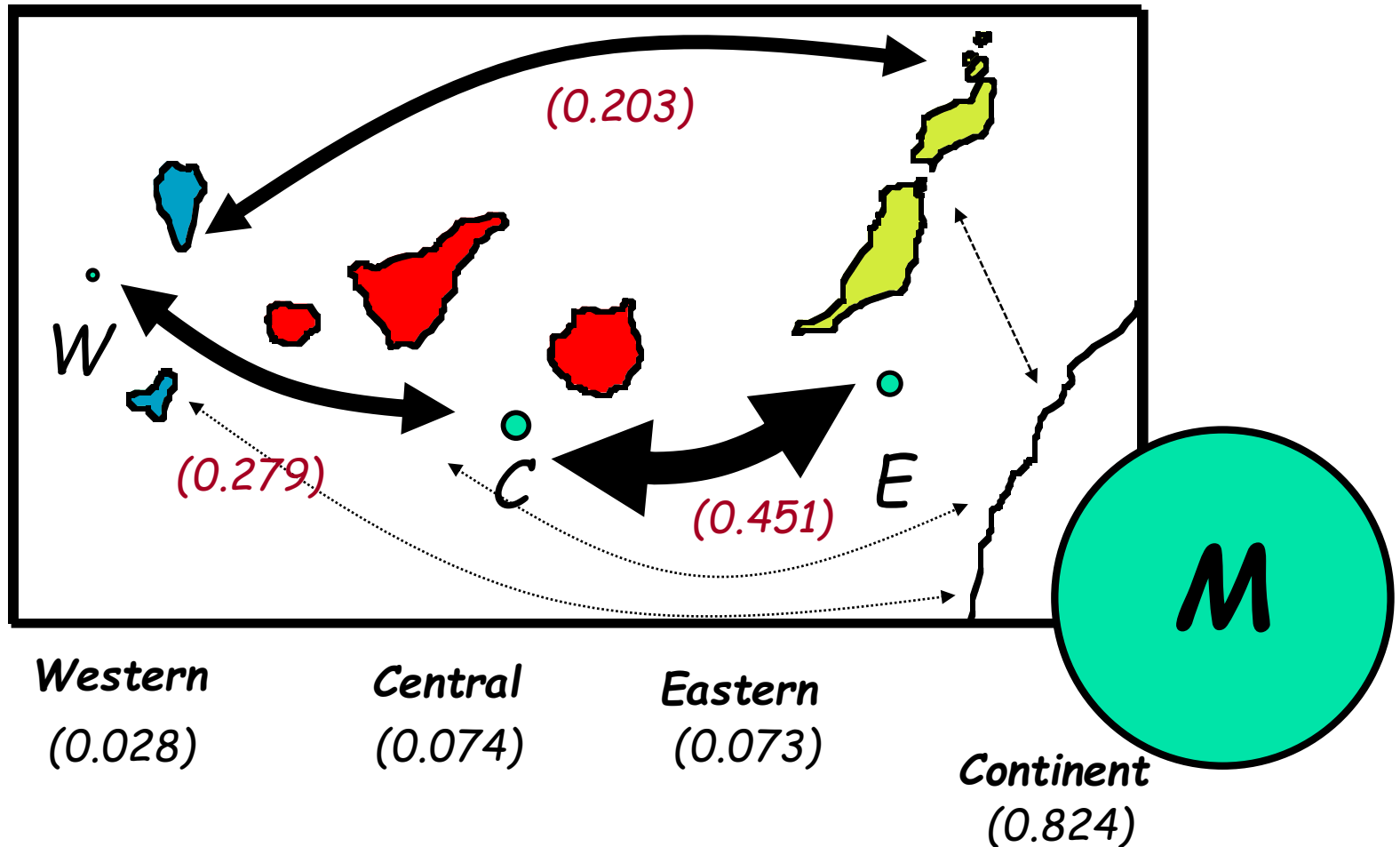
*19 groups, 578 lineages*



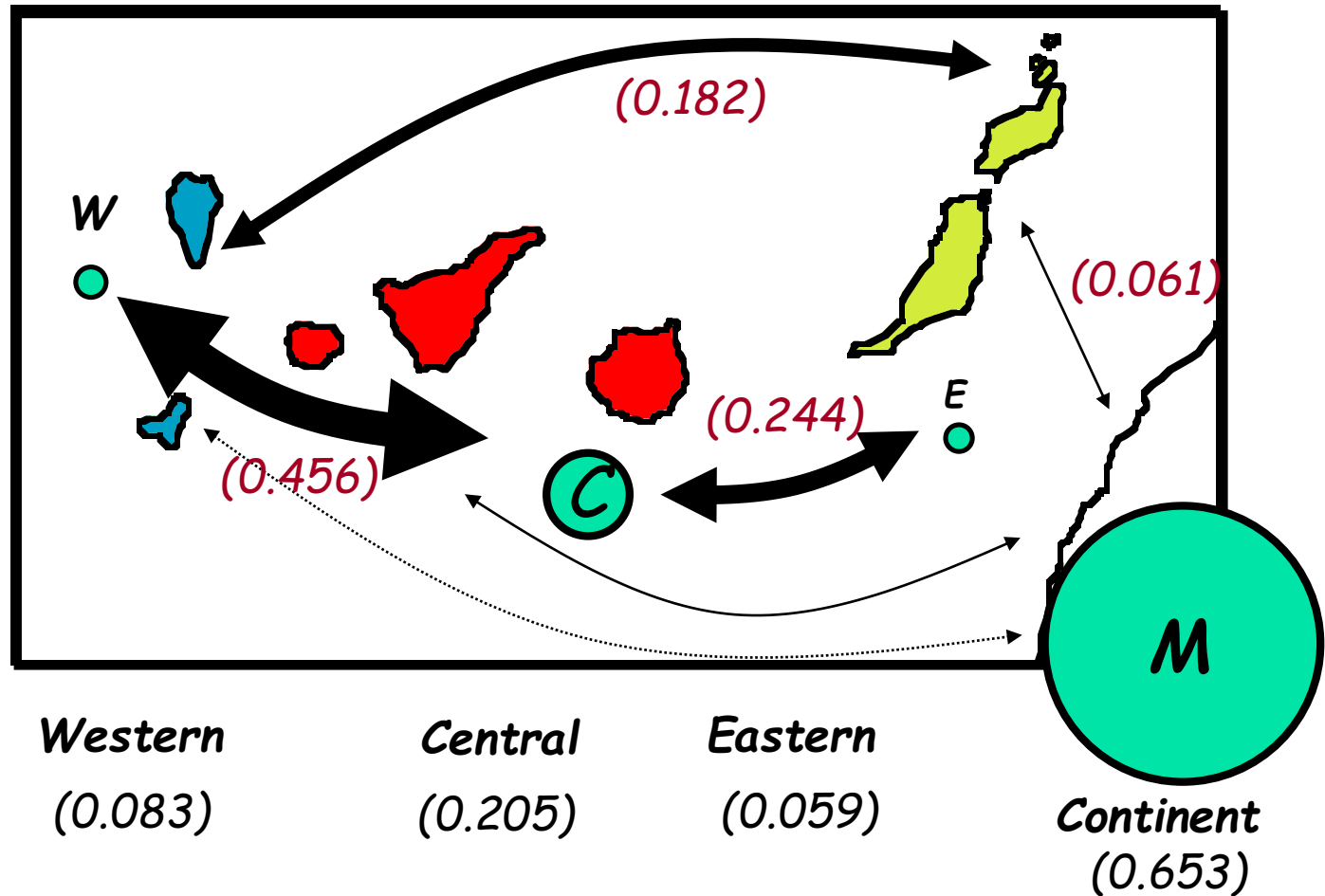
*Javier Fuertes*



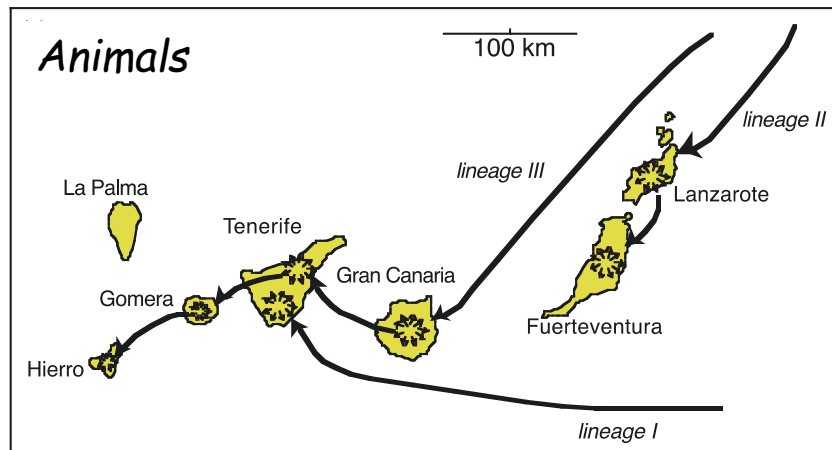
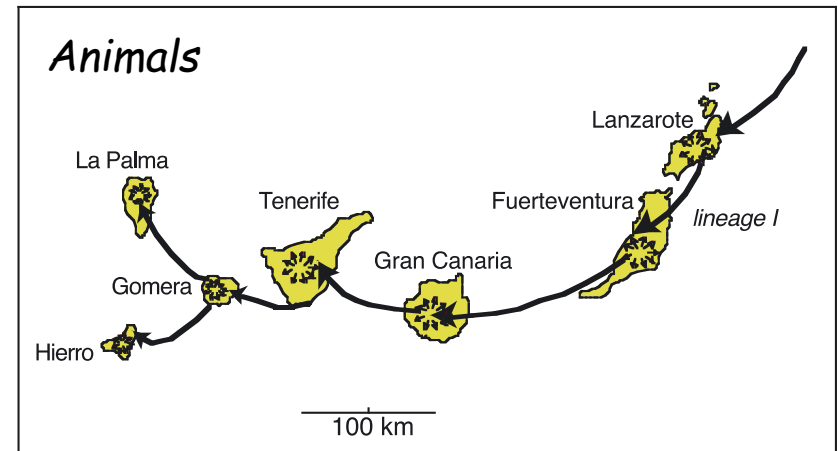
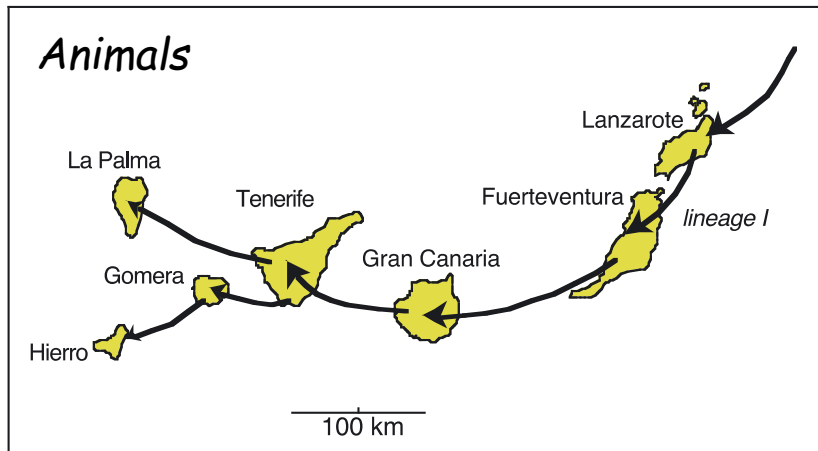
# Colonization patterns (plants)



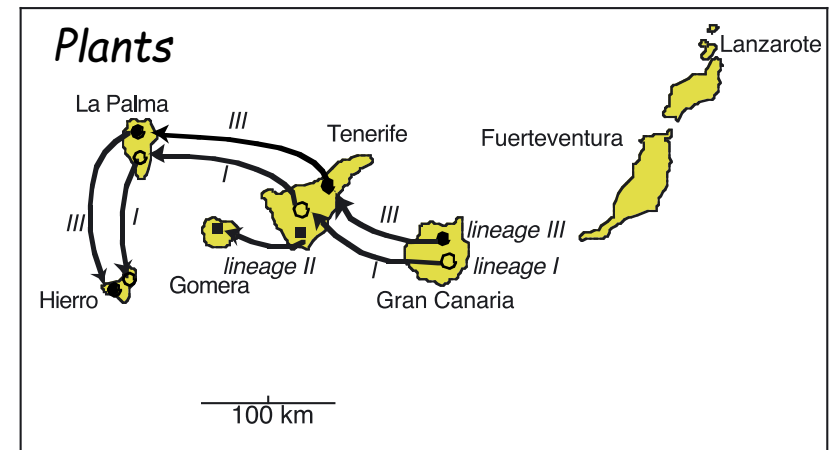
# Colonization patterns (animals)



# Colonization models



*Within-island diversification*  
*Niche evolution*



*Inter-island colonization of*  
*similar biomes (niche conservatism)*

*I move, you change*



# Ecological zones

- Coastal belt
- Open habitat
- Thermophilous forest
- Laurisilva
- Pine forest
- Sub-alpine

## *Ten island-habitat types*

<i>M1</i>	<i>Other Mainland</i>
<i>E2</i>	<i>Eastern-Open</i>
<i>C2</i>	<i>Central-Open</i>
<i>W2</i>	<i>Western Open</i>
<i>C3</i>	<i>Central-laurel forest</i>
<i>W3</i>	<i>Western-laurel forest</i>
<i>C4</i>	<i>Central-pine forest</i>
<i>W4</i>	<i>Western-pine forest</i>
<i>C5</i>	<i>Central-alpine vegetation</i>
<i>W5</i>	<i>Western-alpine vegetation</i>





## *Separating island-hopping and niche-shift rates*

$$r \left\{ \begin{array}{ll} r_i & \text{Shift between islands} \\ r_e & \text{Shift between niches} \\ r_i r_e & \text{Shift between islands and niches} \end{array} \right.$$

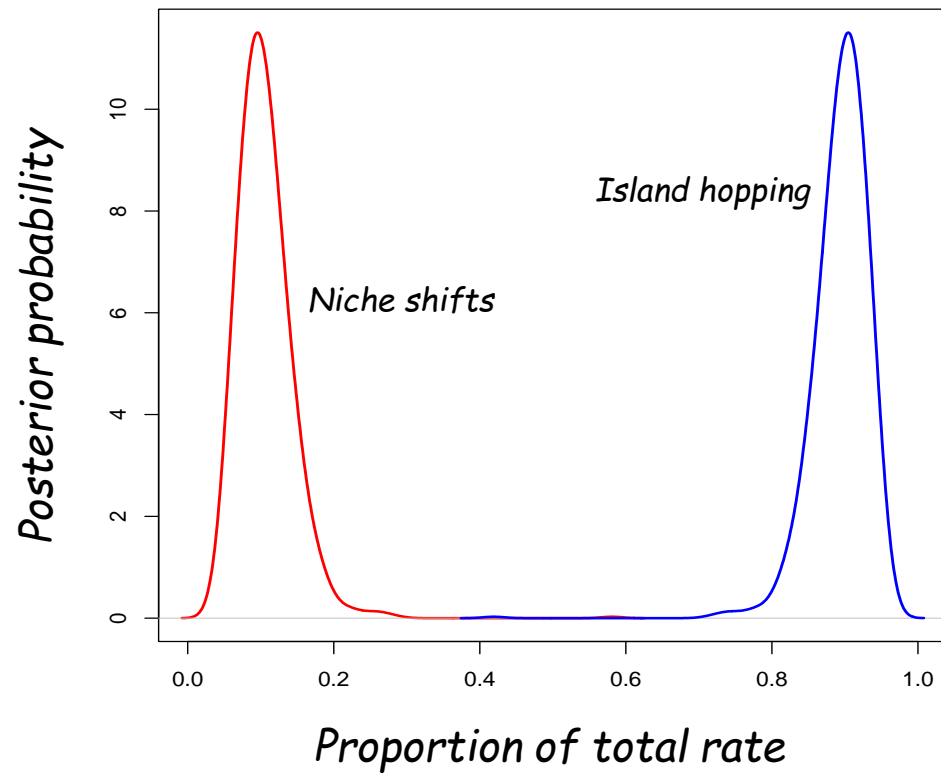
*Standard biogeography model:*

$$r \sim \text{dirichlet}(1, 1, 1, \dots)$$

*Islands-ecology model:*

$$\mu \sim \text{dirichlet}(1, 1)$$

$$r := \text{simplex}(\mu[1], \mu[2], \mu[1] * \mu[2], \dots)$$



## *Separating area and ecology contributions to carrying capacity*

$a, e$                       *Area and ecology components*

$\beta a + (1-\beta)e$       *Linear mix*

$a^\beta e^{(1-\beta)}$               *Power mix*

$S = c A^z$                $z = 0.25 \quad (0.22, 0.28)$

*Standard model:*

$\pi \sim \text{dirichlet}(1, 1, 1, \dots)$

*Power-mix model:*

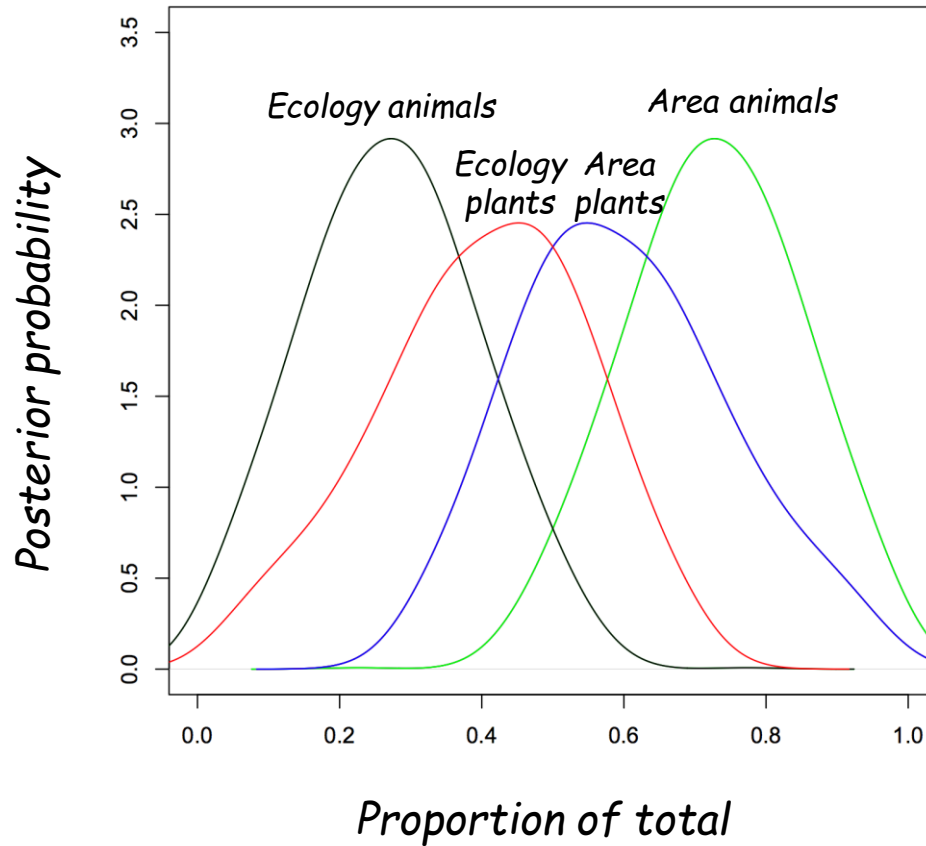
$\text{areaK} \leftarrow \text{simplex}(A_1^z, A_2^z, \dots)$

$\text{ecoK} \sim \text{dirichlet}(1, 1, \dots)$

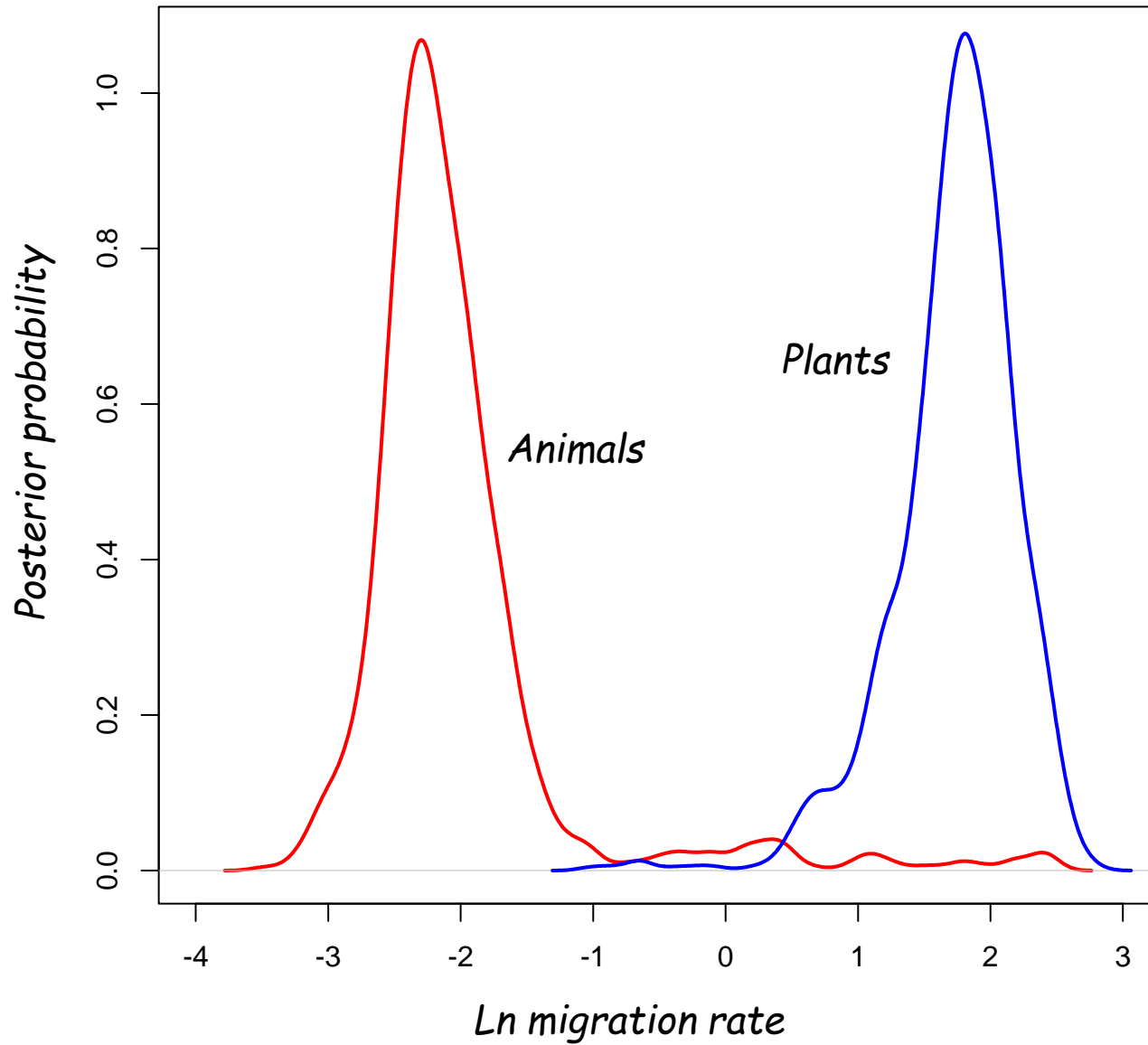
$\text{prop} \sim \text{dirichlet}(1, 1)$

$\pi := \text{powermix}(\text{areaK}, \text{ecoK}, \text{prop})$

*Carrying capacity factored into area  
and other factors (ecology etc)*



## *Comparison of migration rate in animals and plants*



# Summary

- Carrying capacities and relative magnitudes of biotic exchange between islands are similar in animals and plants
- Nevertheless, animal dispersal between islands is an order of magnitude slower than plant dispersal
- Area effects are more important, relatively speaking, in determining island carrying capacities in animals
- Plants shift between islands much more readily than they adapt to new niches
- Phylogenetic graphical models are good tools in modeling and analyzing these phenomena...