

# Heterogeneity in evolutionary processes : structured processes

Joëlle Barido-Sottani

# Tree models in Bayesian inference



$$P(\text{Molecular alignment, Substitution model, Clock model} \mid \text{Time tree, Tree model} \mid \text{ACAC... TCAC... ACAG...}) =$$

Posterior

Likelihood

Probability of  
the tree model

Priors

$$P(\text{ACAC... TCAC... ACAG...} \mid \text{Molecular alignment, Substitution model, Clock model}) \cdot P(\text{Time tree} \mid \text{Tree model}) \cdot P(\text{Molecular alignment, Substitution model, Clock model} \mid \text{Time tree, Tree model})$$

$$P(\text{ACAC... TCAC... ACAG...})$$

ACAC...  
TCAC...  
ACAG...

Molecular alignment



Substitution model



Clock model



Time tree

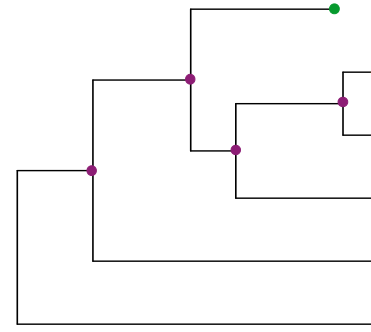


Tree model

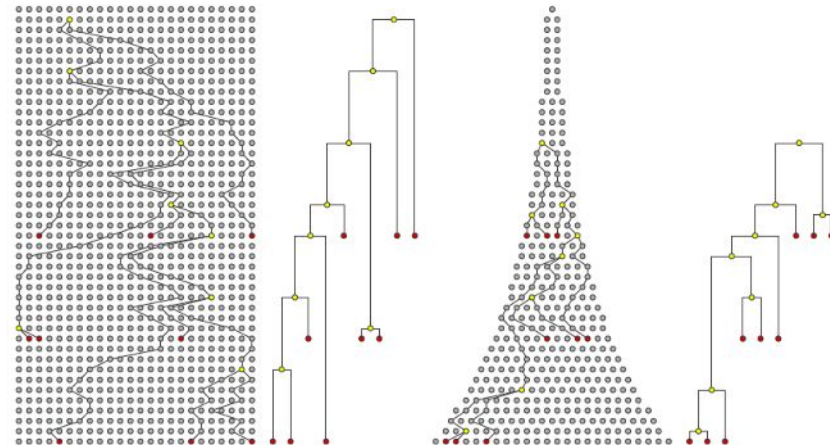


# Phylogenetic models

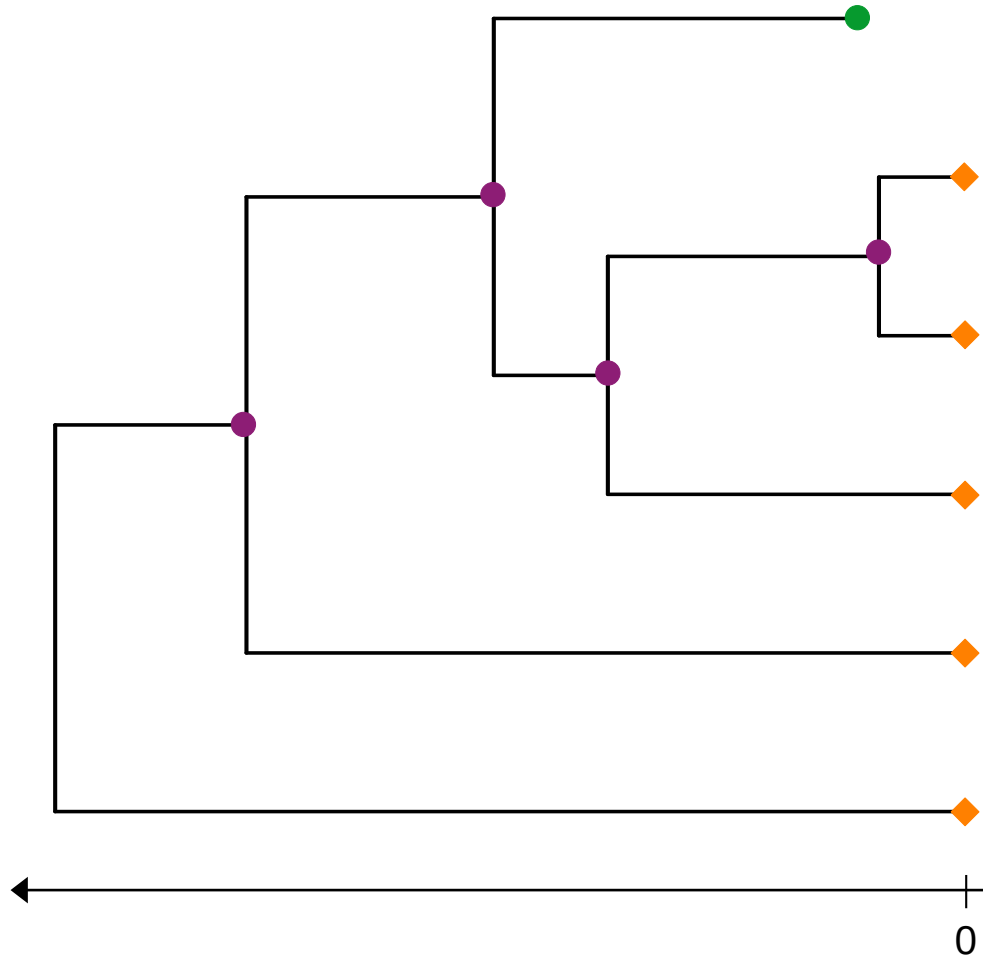
- Birth-death models



- Coalescent models



# Simple birth-death process



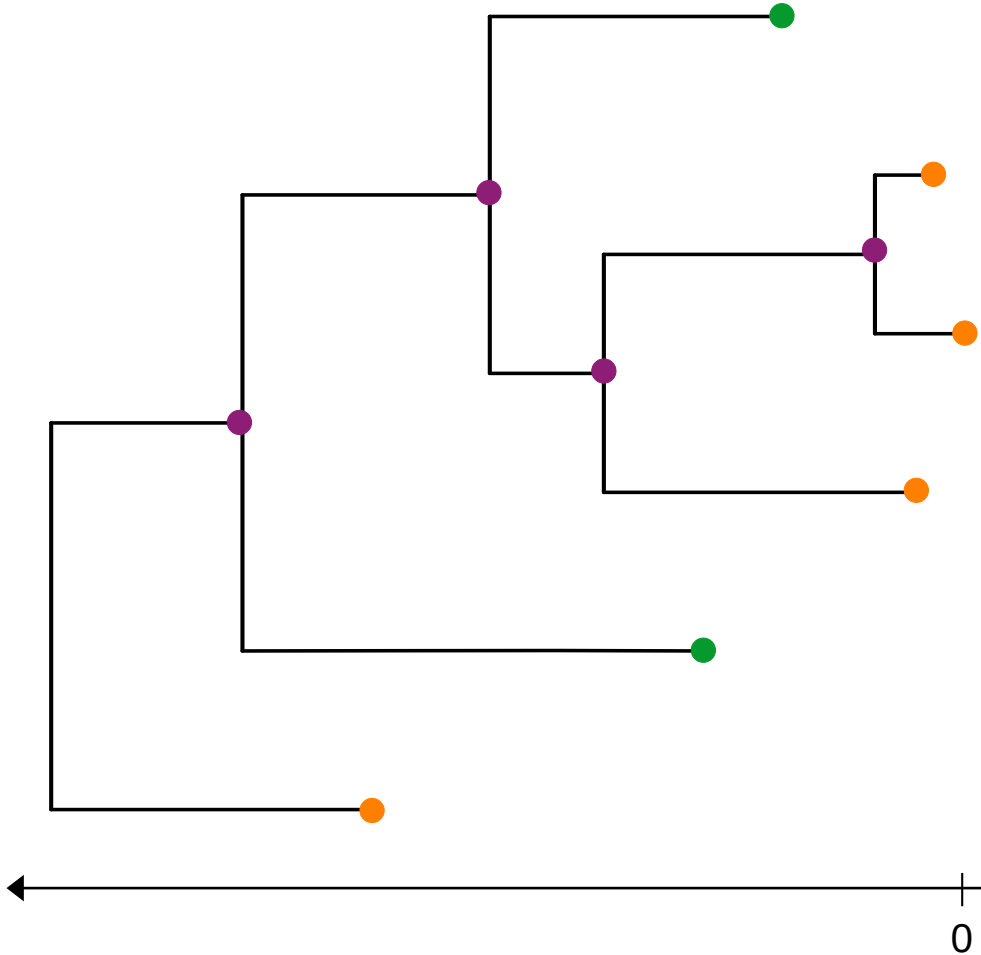
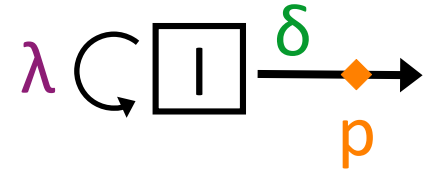
Parameters:

●  $\lambda$  — birth rate (= new lineage appearing)

●  $\mu$  — death rate (= lineage disappearing)

◆  $\rho$  — extant species sampling probability

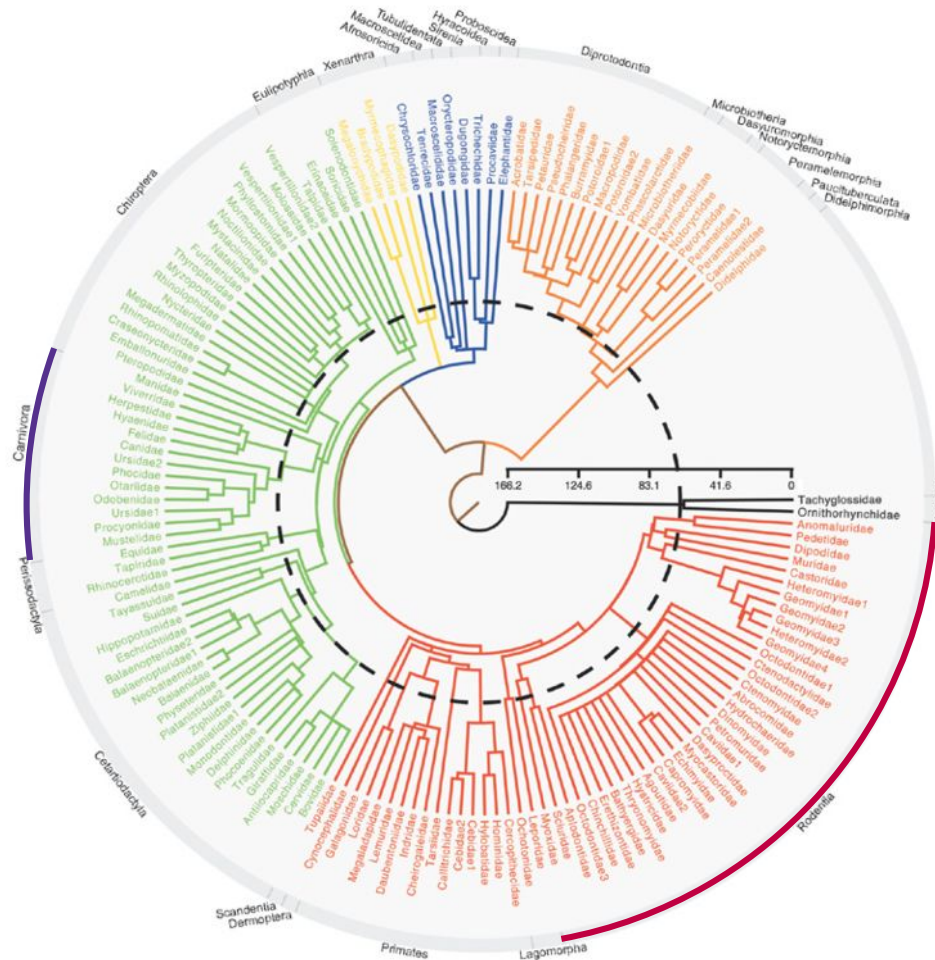
# Birth-death for epidemiology



## Processes:

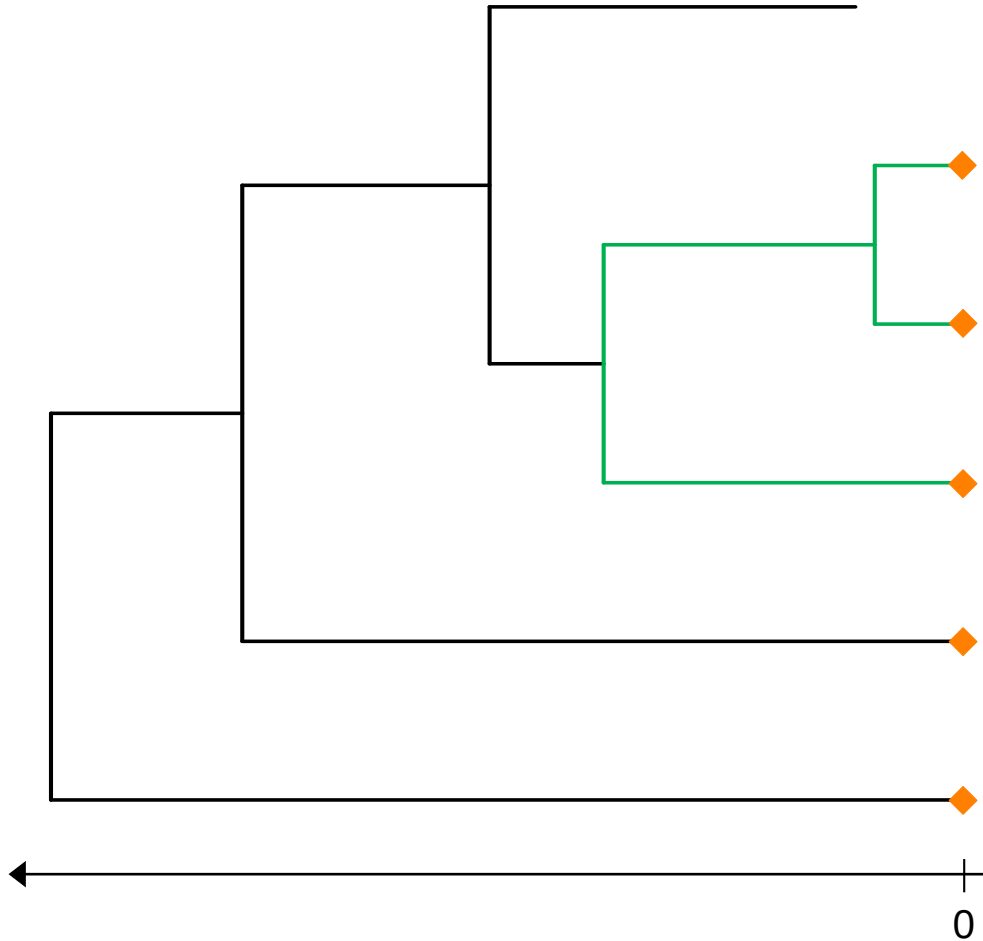
- $\lambda$  — transmission rate
- $\mu = \delta(1-p)$  — rate of recovery without sampling
- $\psi = \delta p$  — rate of recovery with sampling

# Heterogeneity in evolution



- Size discrepancies are evidence of variations in evolutionary processes
- Many traits are proposed to drive variation:
  - body size, mating system, environment, etc.
  - host location, pathogen strain, host behaviour, etc.

# Multi-type birth-death (MTBD) process



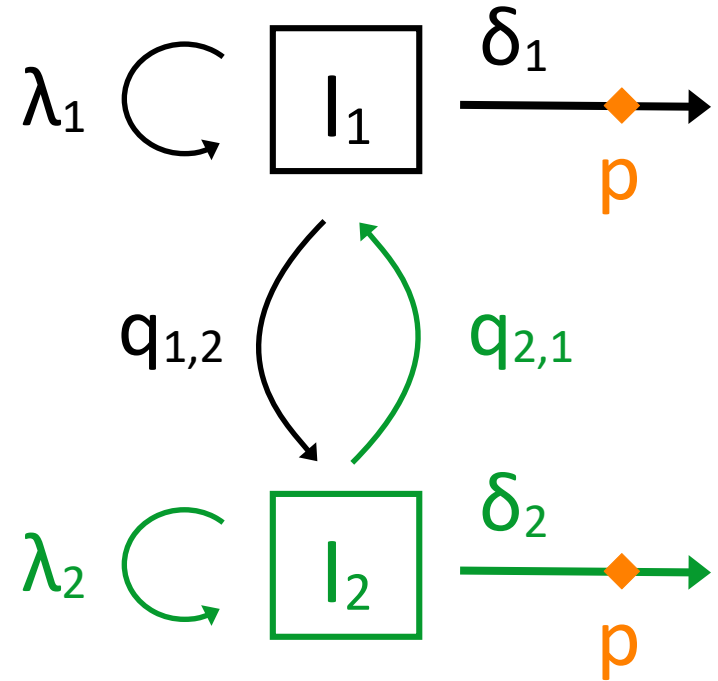
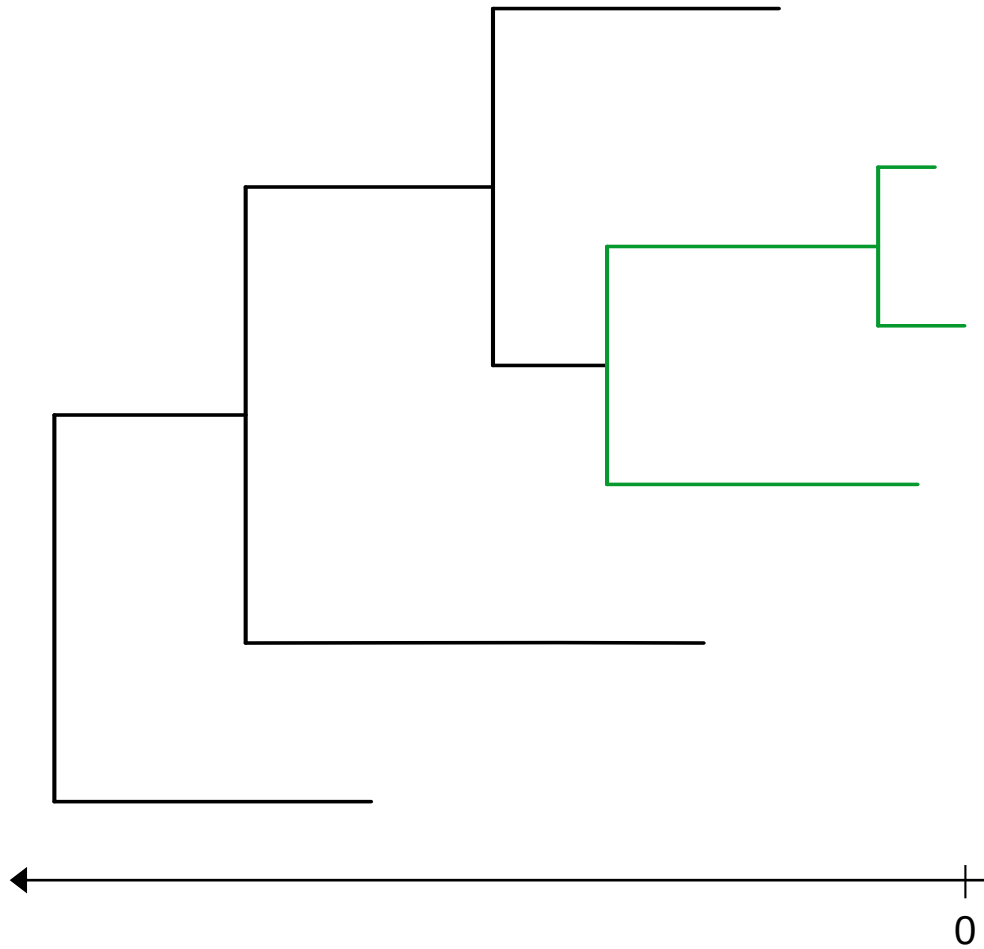
2 types, type 1 & **type 2**

$\lambda_1$  &  **$\lambda_2$**  — birth rates

$\mu_1$  &  **$\mu_2$**  — death rates

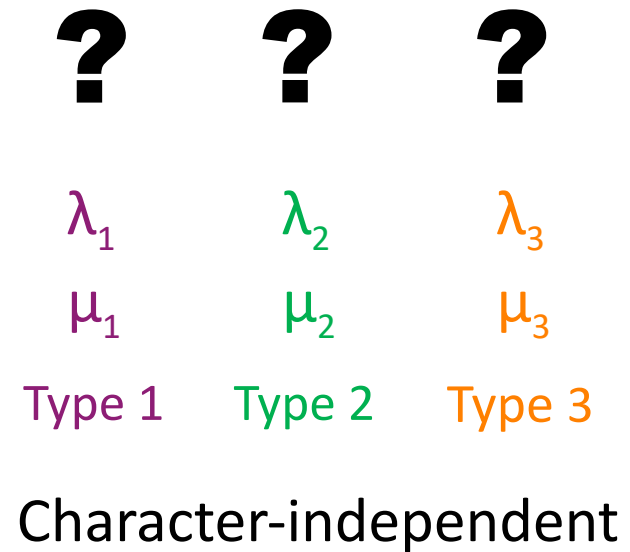
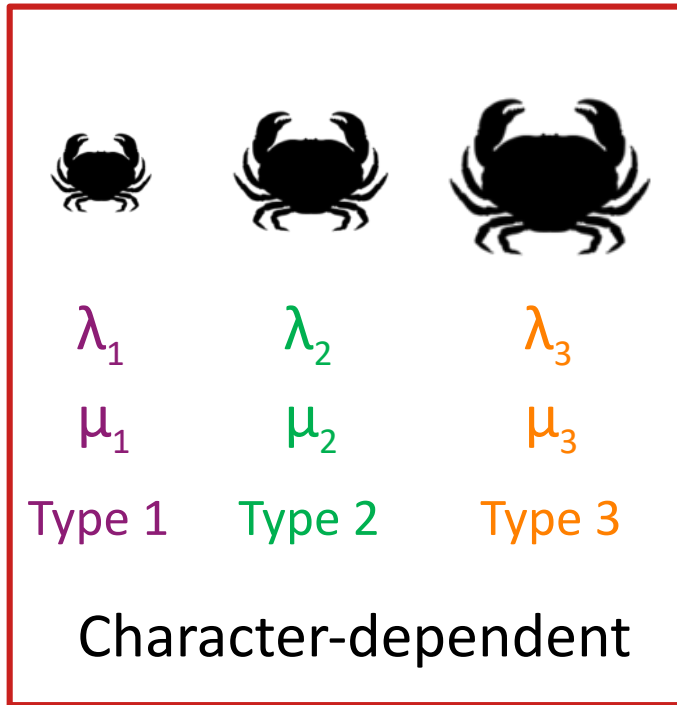
$\rho$  — extant species  
sampling probability

# MTBD process (epidemiology)





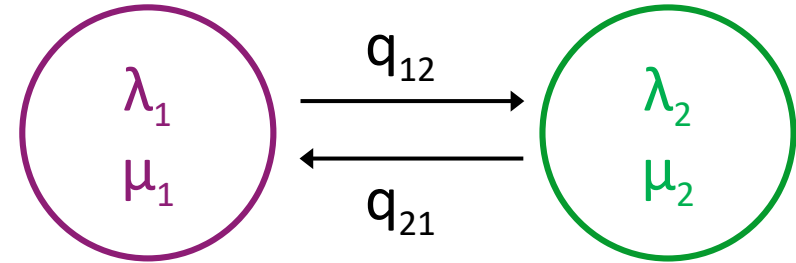
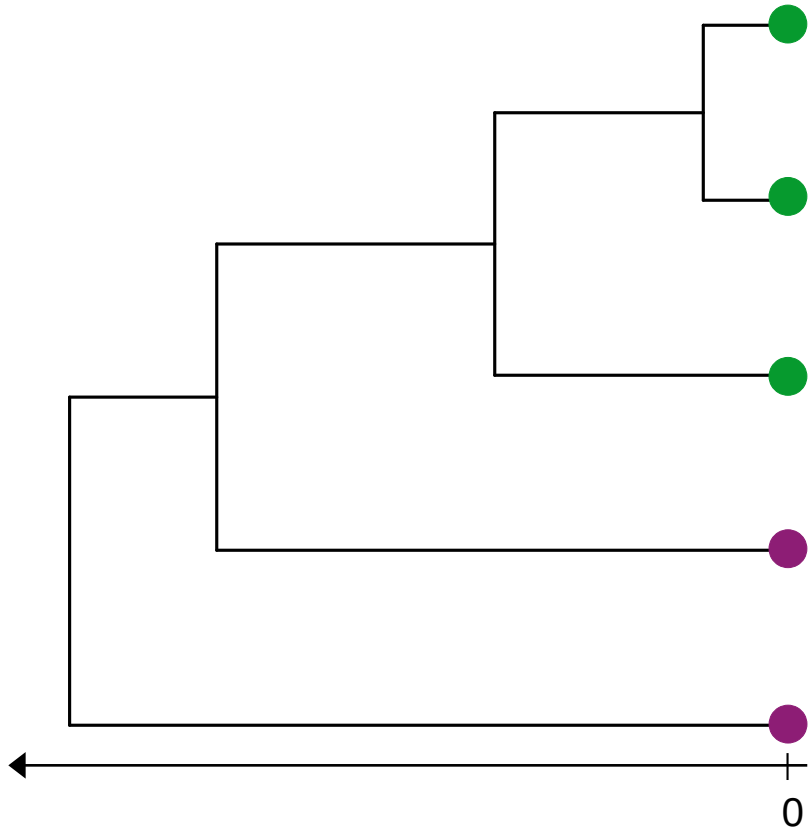
# Character-dependent or independent ?



In a character-dependent model :

- The number of types is known
- The type at the tips is known

# The BiSSE/MuSSE/BDMM model



Parameters of the model:

$\lambda_i$  – birth rates

$\mu_i$  – death rates

$q_{ij}$  – transition rates

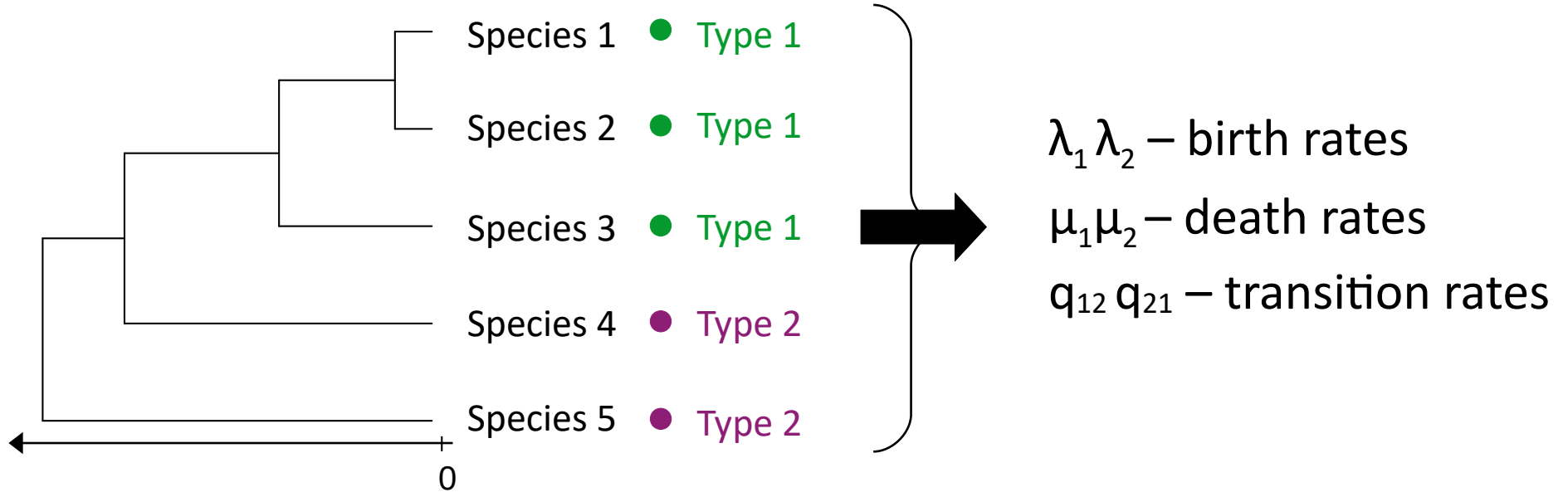
$\rho/p$  – sampling probability

Maddison *et al.* **Sys. Bio.** 2007

Fitzjohn *et al.* **Sys. Bio.** 2009

Kühnert *et al.* **MBE** 2016

# SSE/BDMM inference

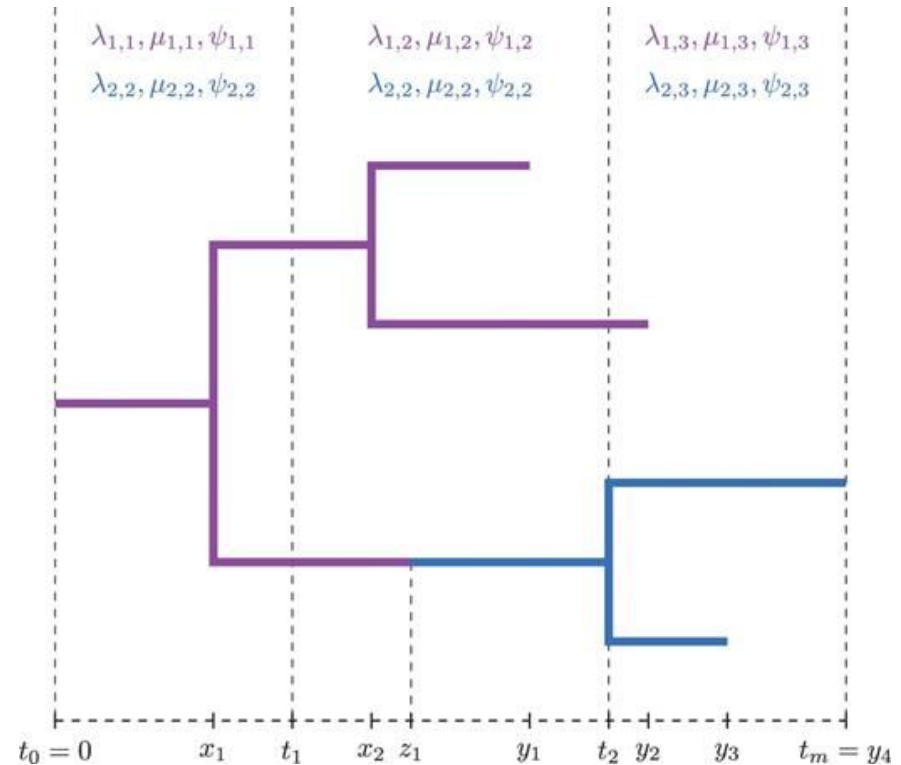


Important assumption: the evolutionary processes in the complete phylogeny (including non-sampled parts) are identical to the processes in the reconstructed phylogeny.

# BDMM extensions



- Integration with the skyline model: piecewise-constant rates per type
- Sampling proportion per type:  $p_i$
- Cross-type birth events:  $\lambda_{i,j}$



# New: BDMM-Prime

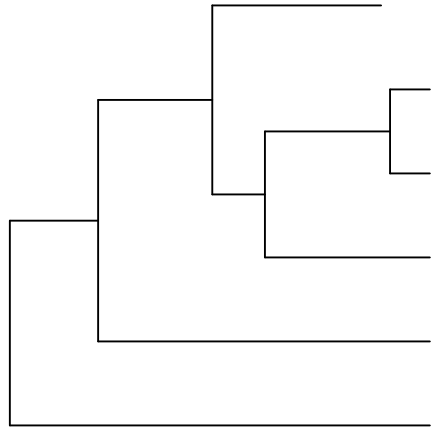


- Can estimate the full (unsampled) LTT
- Uses advanced computational techniques for faster performance
- Can estimate the event history or integrate over it
- New improved BEAUti interface for setup

## Bayesian Phylodynamic Inference of Multitype Population Trajectories Using Genomic Data

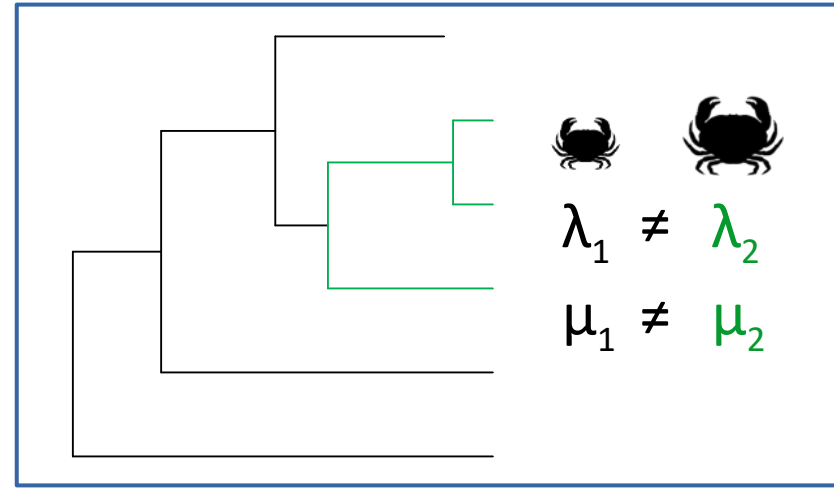
Timothy G. Vaughan <sup>1,2,\*</sup> Tanja Stadler <sup>1,2</sup>

# Model selection issues



$$\lambda_1 = \lambda_2$$

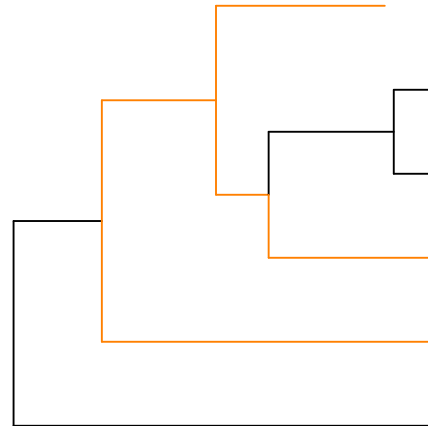
$$\mu_1 = \mu_2$$



$$\lambda_1 \neq \lambda_2$$

$$\mu_1 \neq \mu_2$$

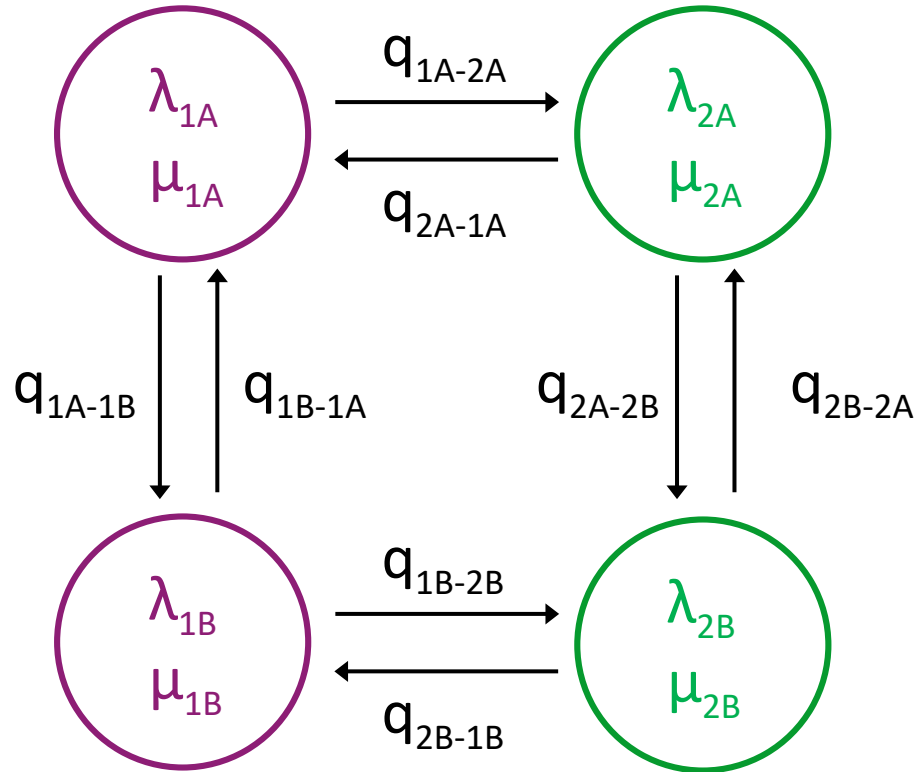
Rabosky & Goldberg 2015, Sys. Bio.



$$\lambda_1 \neq \lambda_2$$

$$\mu_1 \neq \mu_2$$

# The HiSSE model



- Hybrid model with a hidden character (A/B) added to the user-chosen trait (1/2)
- Allows to distinguish whether the user-chosen character is linked to the rate variation
- Only single transitions are allowed (no diagonal)
- Remaining issues:
  - The number of values for the hidden character is chosen by the user
  - Higher complexity of the model

# Examples: character-driven diversification

## RESEARCH ARTICLE



### Evidence linking life-form to a major shift in diversification rate in *Crassula*

Meng Lu<sup>1,2</sup> | Marc Fradera-Soler<sup>1,3</sup> | Félix Forest<sup>1</sup> |  
Timothy G. Barraclough<sup>2,4</sup> | Olwen M. Grace<sup>1</sup>

## ORIGINAL ARTICLE

Ecological  
Entomology

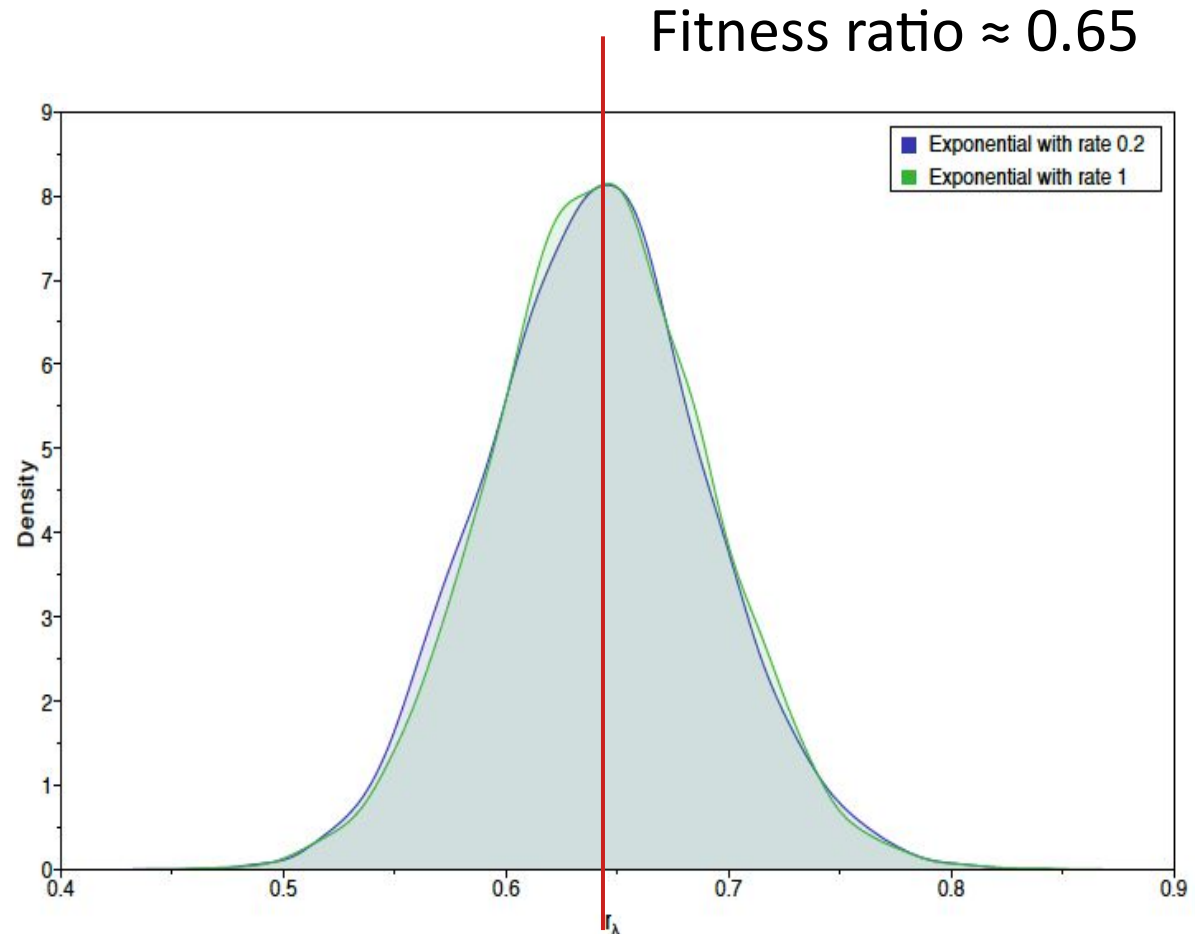
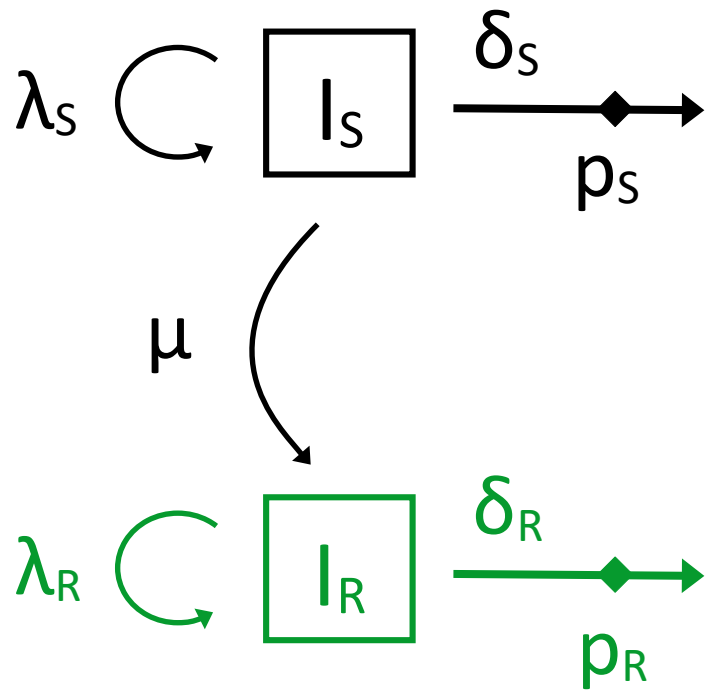


### Fast species diversification among dragonflies (Anisoptera: Odonata: Insecta) inhabiting lentic environments regardless of wing pigmentation

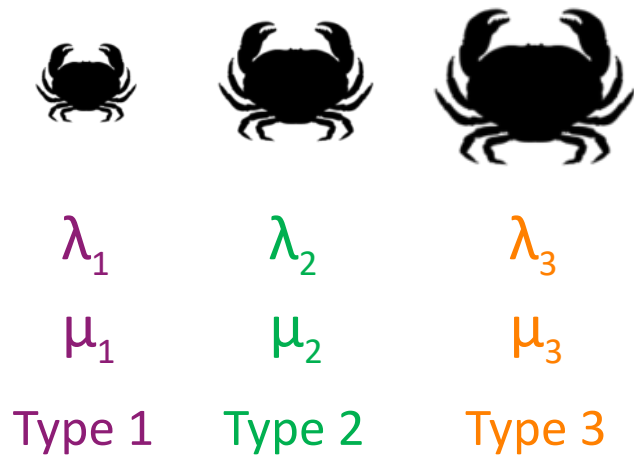
Benjamín Padilla-Morales<sup>1</sup> | Paola Cornejo-Páramo<sup>1</sup> | Oscar García-Miranda<sup>2</sup> |  
Aldo Issac Carrillo Muñoz<sup>2</sup> | Andrea Nieto López<sup>2</sup> | Daniel L. Castillo-Morales<sup>3</sup> |  
Gustavo Wappler Barragán<sup>2</sup> | Araxi O. Urrutia<sup>1,3</sup> | Martín Alejandro Serrano-Meneses<sup>2</sup>



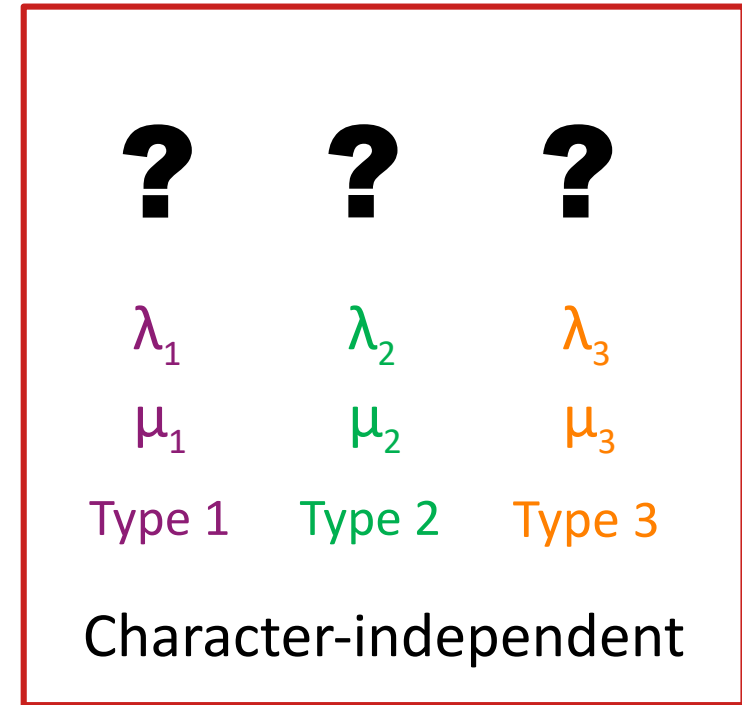
# Example: fitness of resistant tuberculosis



# Character-dependent or independent ?



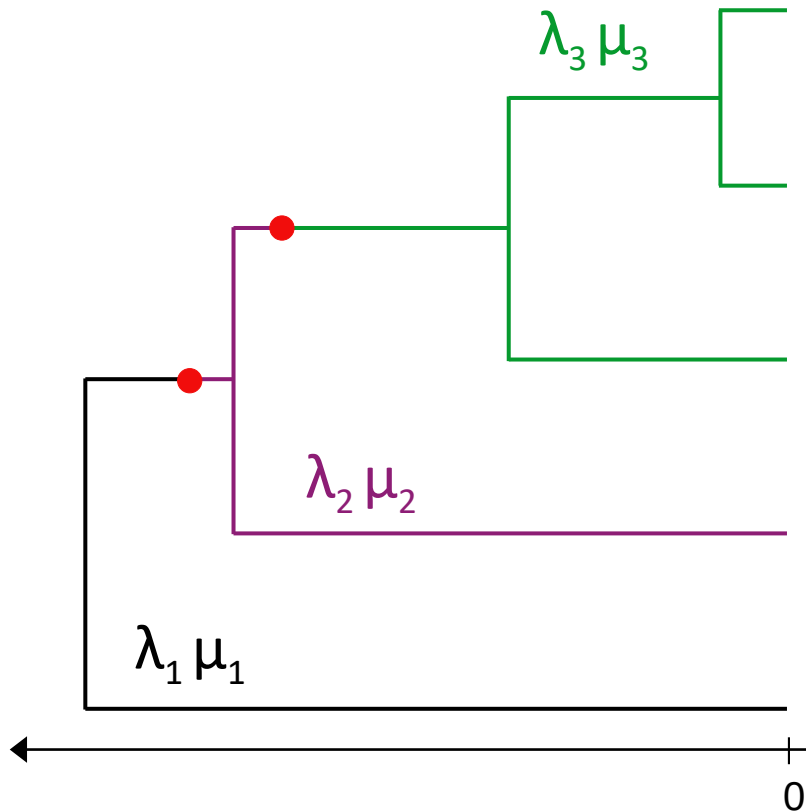
Character-dependent



In a character-dependent model :

- The number of types is known
- The type at the tips is known

# BAMM/MSBD model



- Character-independent version of SSE
- New estimated parameters:
  - N total number of types
  - Types of edges and tips
- Simplified transition process: ●
  - Each transition is a new type (BAMM)
  - Constant transition rate  $\gamma$  (MSBD)
- Assumes that all types appear in the sampled tree – no unseen types

Rabosky *et al.* **Nat. Comm.** 2013  
Barido-Sottani *et al.* **Sys. Bio.** 2020

# Simplifying the model

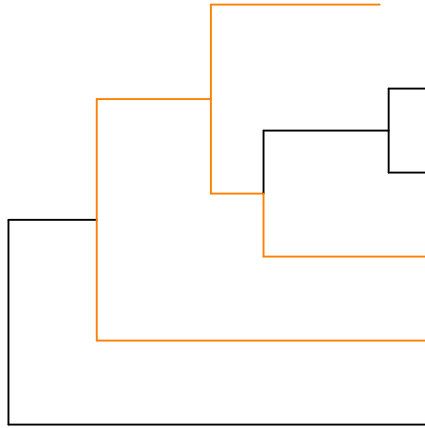
Character-dependent



$$\begin{array}{ccc} \lambda_1 & \approx & \lambda_2 & \lambda_3 \\ \mu_1 & \approx & \mu_2 & \mu_3 \end{array}$$



$$\begin{array}{cc} \lambda_1 & \lambda_2 \\ \mu_1 & \mu_2 \end{array}$$



Character-independent



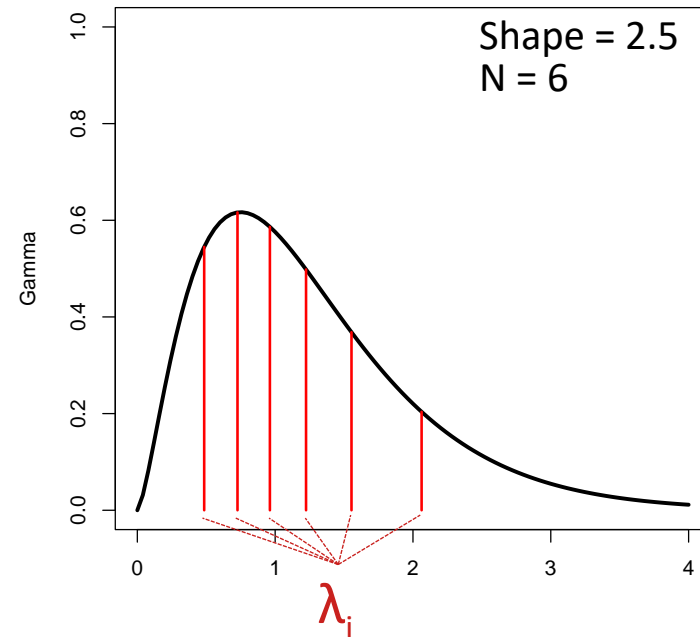
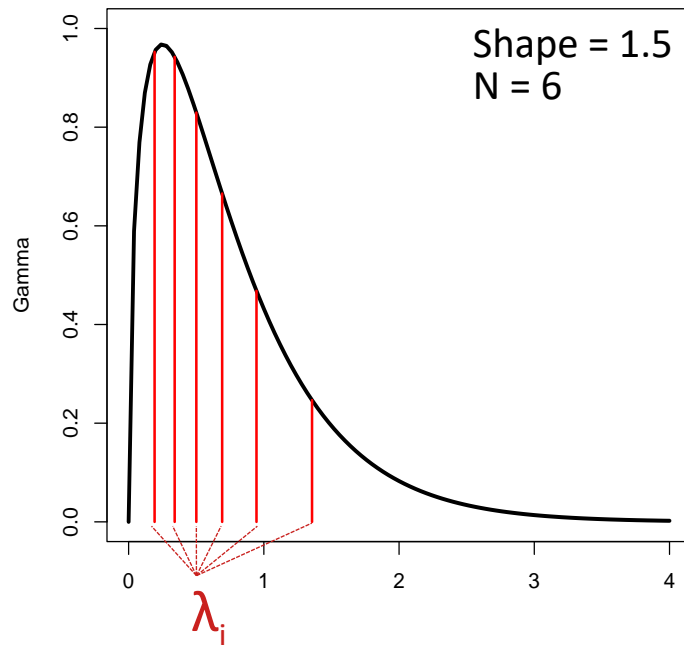
~~$$\begin{array}{ccc} \lambda_1 & \approx & \lambda_2 & \lambda_3 \\ \mu_1 & \approx & \mu_2 & \mu_3 \end{array}$$~~



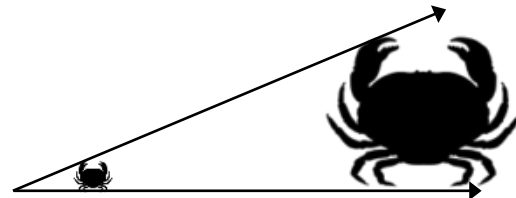
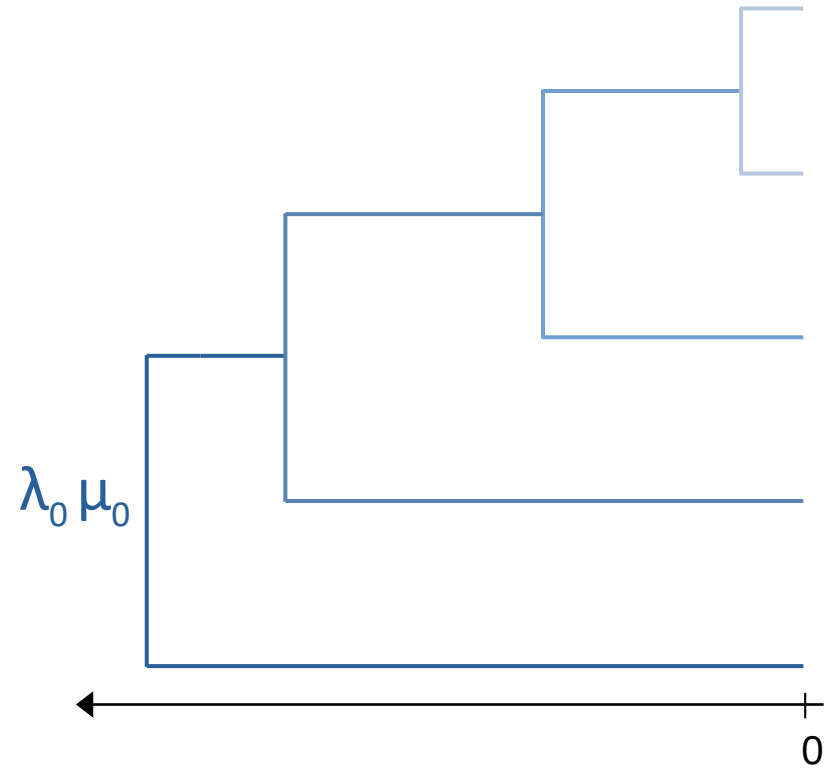
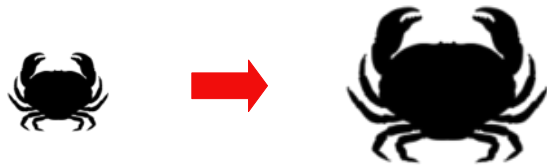
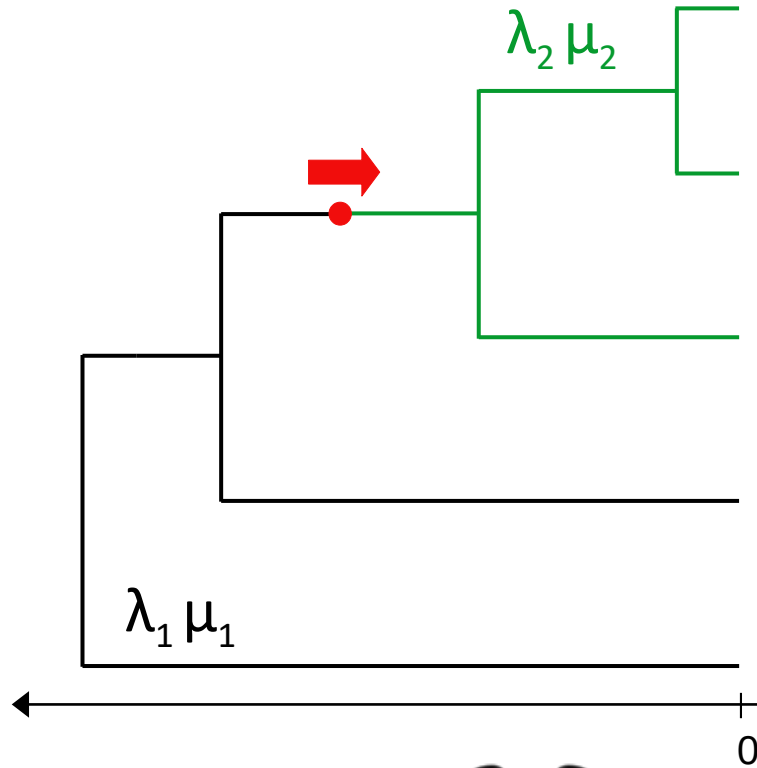
$$\begin{array}{cc} \lambda_1 & \lambda_2 \\ \mu_1 & \mu_2 \end{array}$$

# RevBayes model

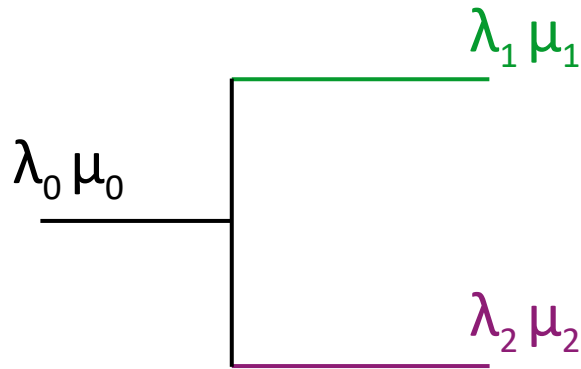
- Ordered types based on a Gamma distribution
- Fixed number of types N
- Simplified model: rates are not estimated, but determined by the shape of the Gamma distribution



# Going beyond types



# ClaDS model



$$\lambda_1 = \text{LogNormal}(\lambda_0 \times \alpha, \sigma)$$

$$\lambda_2 = \text{LogNormal}(\lambda_0 \times \alpha, \sigma)$$

- Continuous evolution process driven by a lognormal distribution
- New estimated parameters:
  - Initial rates at the root  $\lambda_0$  and  $\mu_0$
  - Lognormal parameters  $\alpha$  and  $\sigma$
  - Birth rates for each edge  $\lambda_i$
- Two parameterizations for  $\mu$ 
  - Lognormal process with  $\alpha_\mu$  and  $\sigma_\mu$
  - Assumption of constant turnover:  
 $\mu_i / \lambda_i = \mu_0 / \lambda_0$

Maliet *et al.* **Nat. Eco. Evo.** 2019

Maliet & Morlon **Sys. Bio.** 2021

Barido-Sottani & Morlon **Sys. Bio.** 2023

# Examples: character-driven diversification

ARTICLE



<https://doi.org/10.1038/s41467-020-16498-w>

OPEN

## Trophic innovations fuel reef fish diversification

Alexandre C. Siqueira<sup>1</sup> <sup>✉</sup>, Renato A. Morais<sup>1,2</sup> , David R. Bellwood<sup>1,2</sup>  & Peter F. Cowman<sup>1</sup> 

Nat. Comm. 2020

## No link between population isolation and speciation rate in squamate reptiles

Sonal Singhal<sup>a,1</sup>, Guarino R. Colli<sup>b</sup> , Maggie R. Grundler<sup>c,d</sup>, Gabriel C. Costa<sup>e</sup> , Ivan Prates<sup>f,g</sup>, and Daniel L. Rabosky<sup>f,g,1</sup>

PNAS 2022



# So – character-dependent or independent ?

## Character-dependent / hybrid

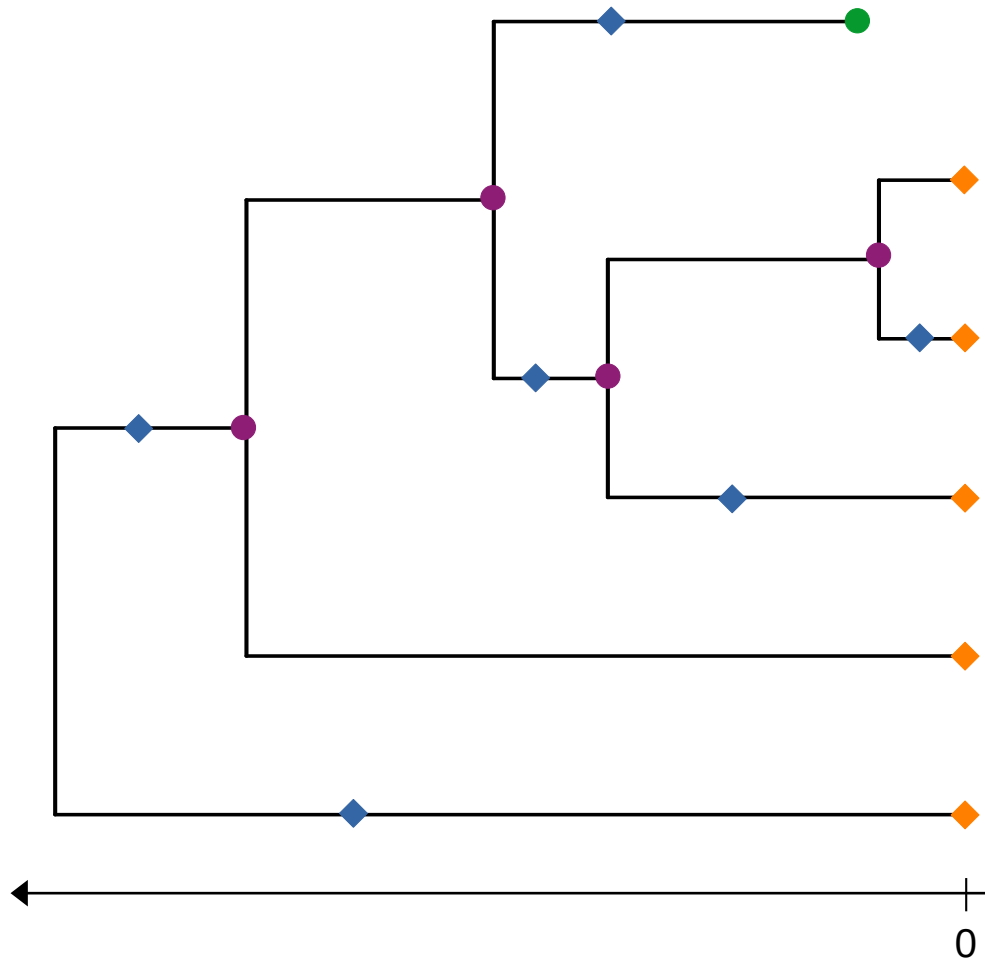
- Allow more complex transition processes
- Are more easily interpreted
- Are very dependent on the choice and accuracy of trait

## Character-independent

- Usually have to make simplifying assumptions
- Do not give direct answers
- Are not constrained by trait information or hypothesis

What is your hypothesis ?  
What are you trying to find out ?

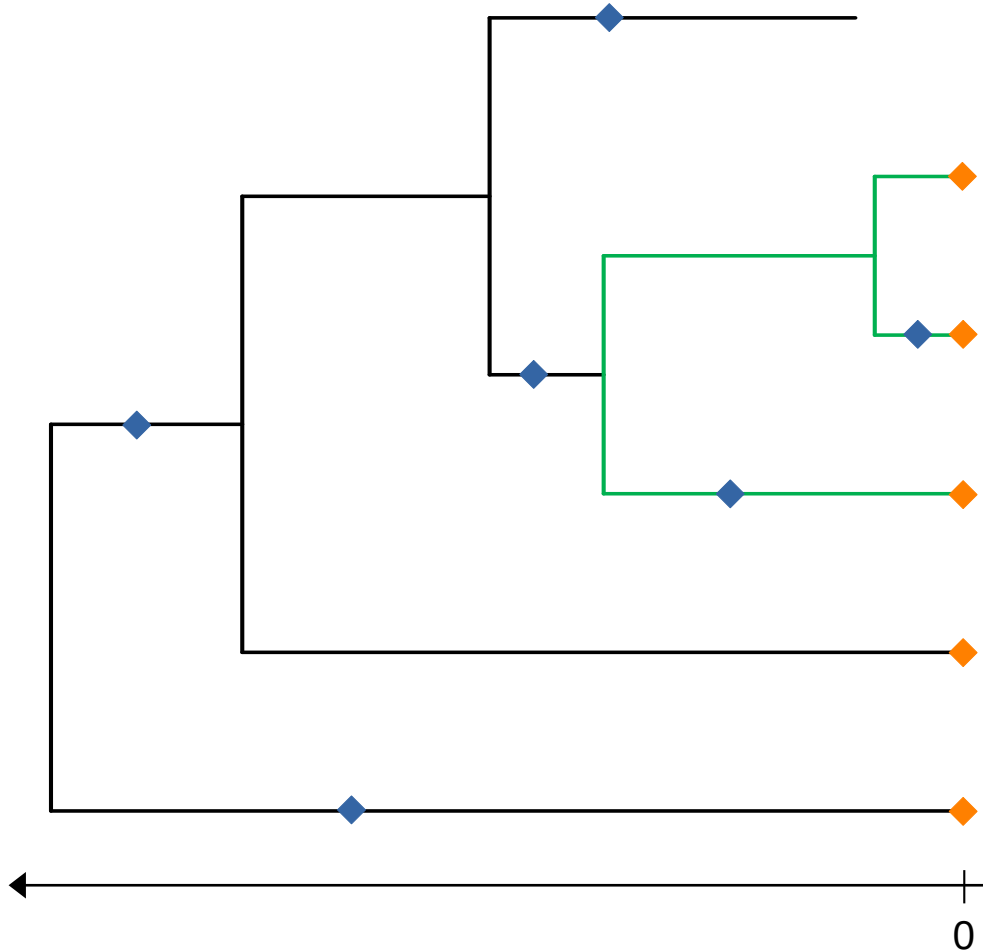
# Integration of fossil/SA data



Parameters:

- $\lambda$  — birth rate
- $\mu$  — death rate
- ◆  $\psi$  — fossilization rate
- ◆  $\rho$  — extant species sampling probability

# Multi-type FBD/SA process



2 types, type 1 & **type 2**

$\lambda_1$  &  $\lambda_2$  — birth rates

$\mu_1$  &  $\mu_2$  — death rates

$\psi_1$  &  $\psi_2$  — fossilization rates

$\rho$  — extant species  
sampling probability

# In summary

- Empirical data supports widespread variation in evolutionary processes, which can be modeled using multi-type birth-death processes
- Multi-type birth-death processes come in two main categories:
  - Character-dependent: uses more information **but** subject to model selection issues
  - Character-independent: more powerful, more expensive and more difficult to interpret
- These models are still a very active area of research and development (extension to continuous processes, integration of fossils, interpretation of results, etc.)

# In summary (BEAST2)

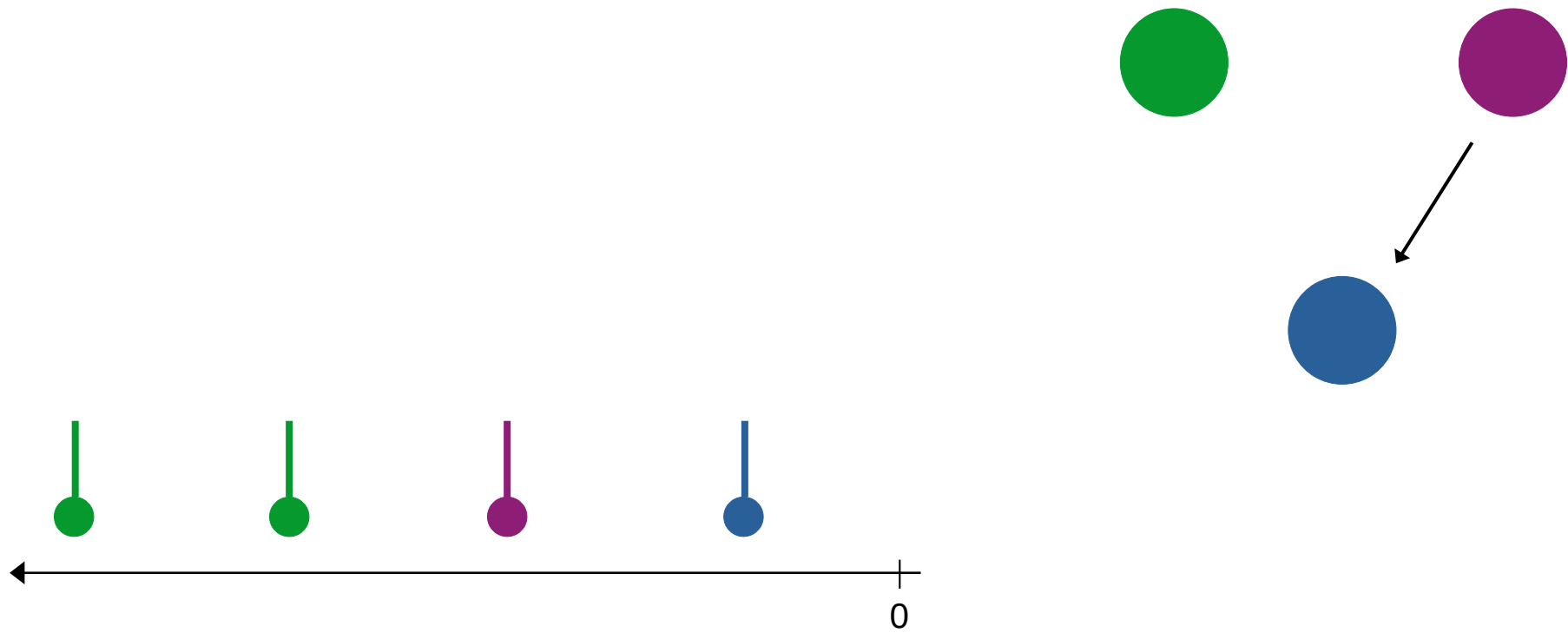


- **Birth-Death-Migration Model (BDMM) package**
  - Character-dependent SSE
  - Includes time-dependent changes (skyline model)
  - Includes sampled ancestors
- **Multi-State Birth-Death (MSBD) package**
  - Character-independent SSE
  - Includes sampled ancestors, starting from v1.3.0
- **Cladogenetic Diversification rate Shift (ClaDS) package**
  - Progressive autocorrelated rate variations
  - Inclusion of sampled ancestors in development

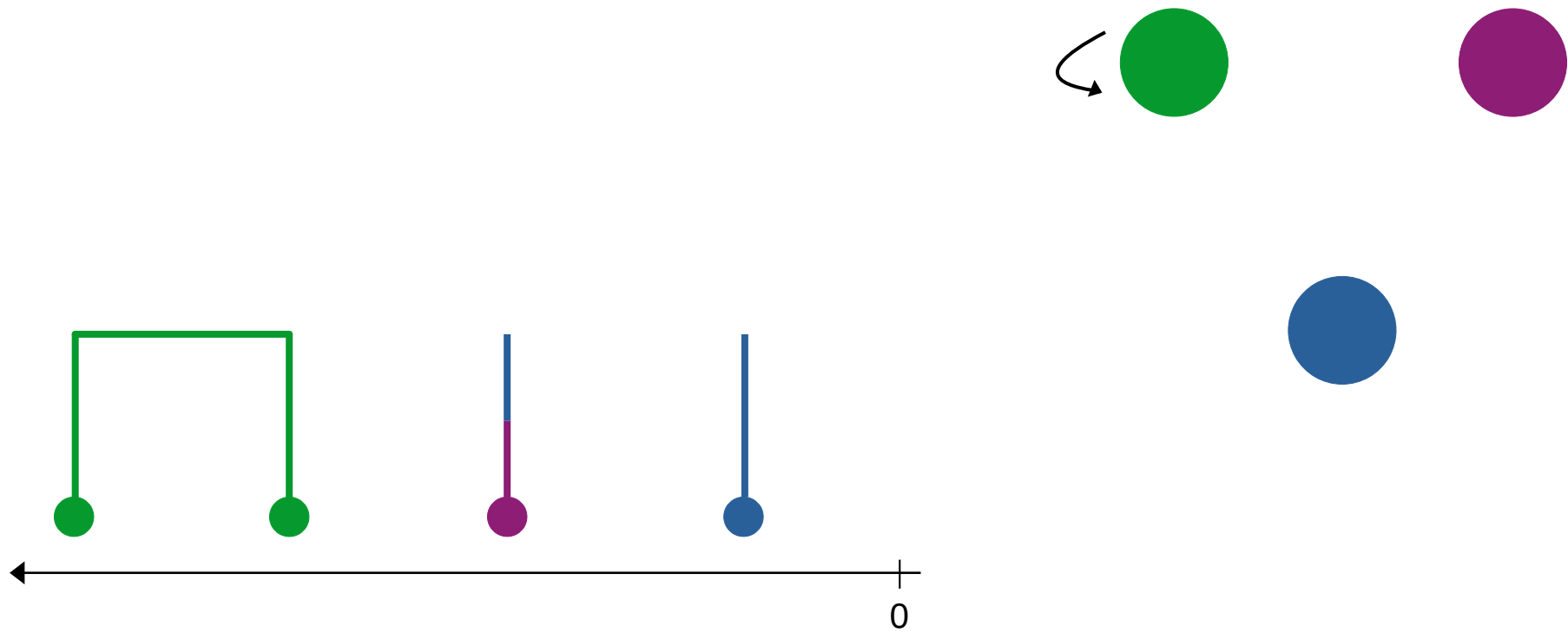
# What about coalescent models?

- Basic coalescent model assumes **exchangeability**: all lineages can coalesce with each other
- Structured coalescent model:
  - $n$  subpopulations with sizes  $N_{e1}, \dots, N_{en}$
  - only lineages of the same subpopulation coalesce
  - adds migration events: one lineage moves from  $i$  to  $j$

# Structured coalescent in action

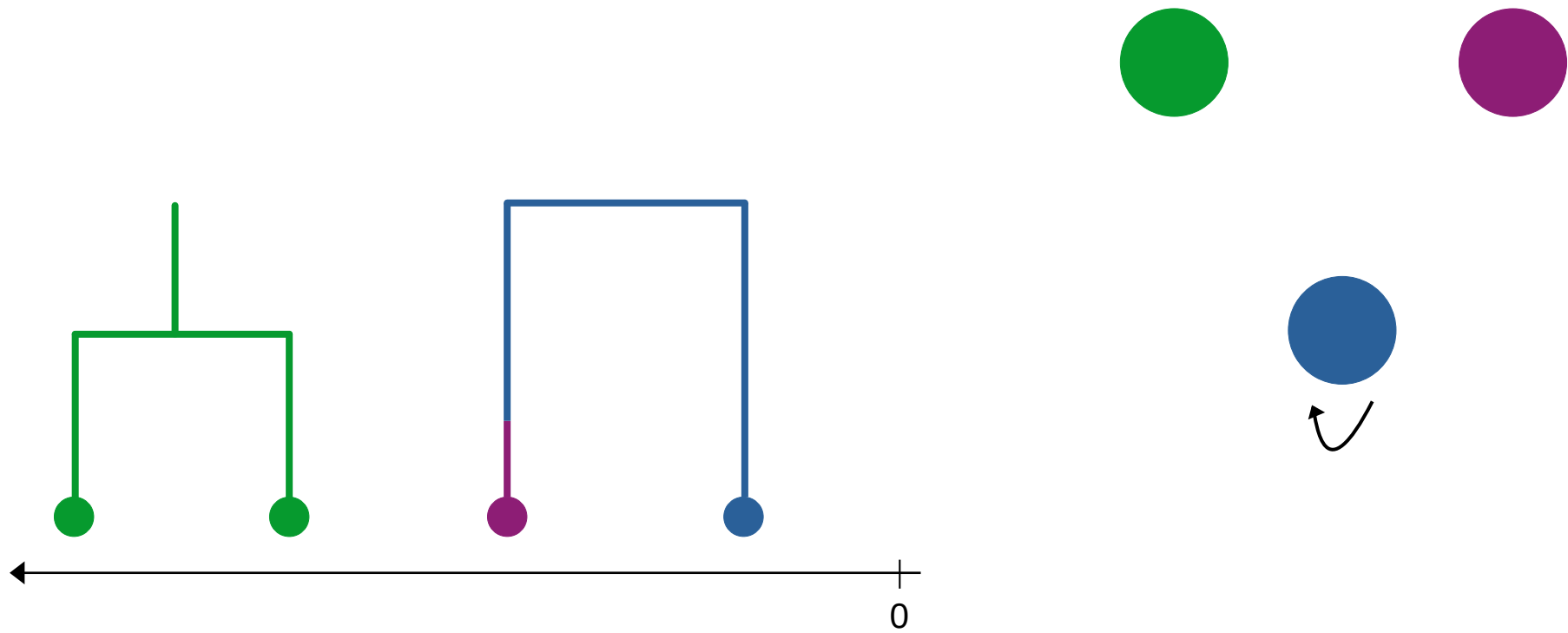


# Structured coalescent in action

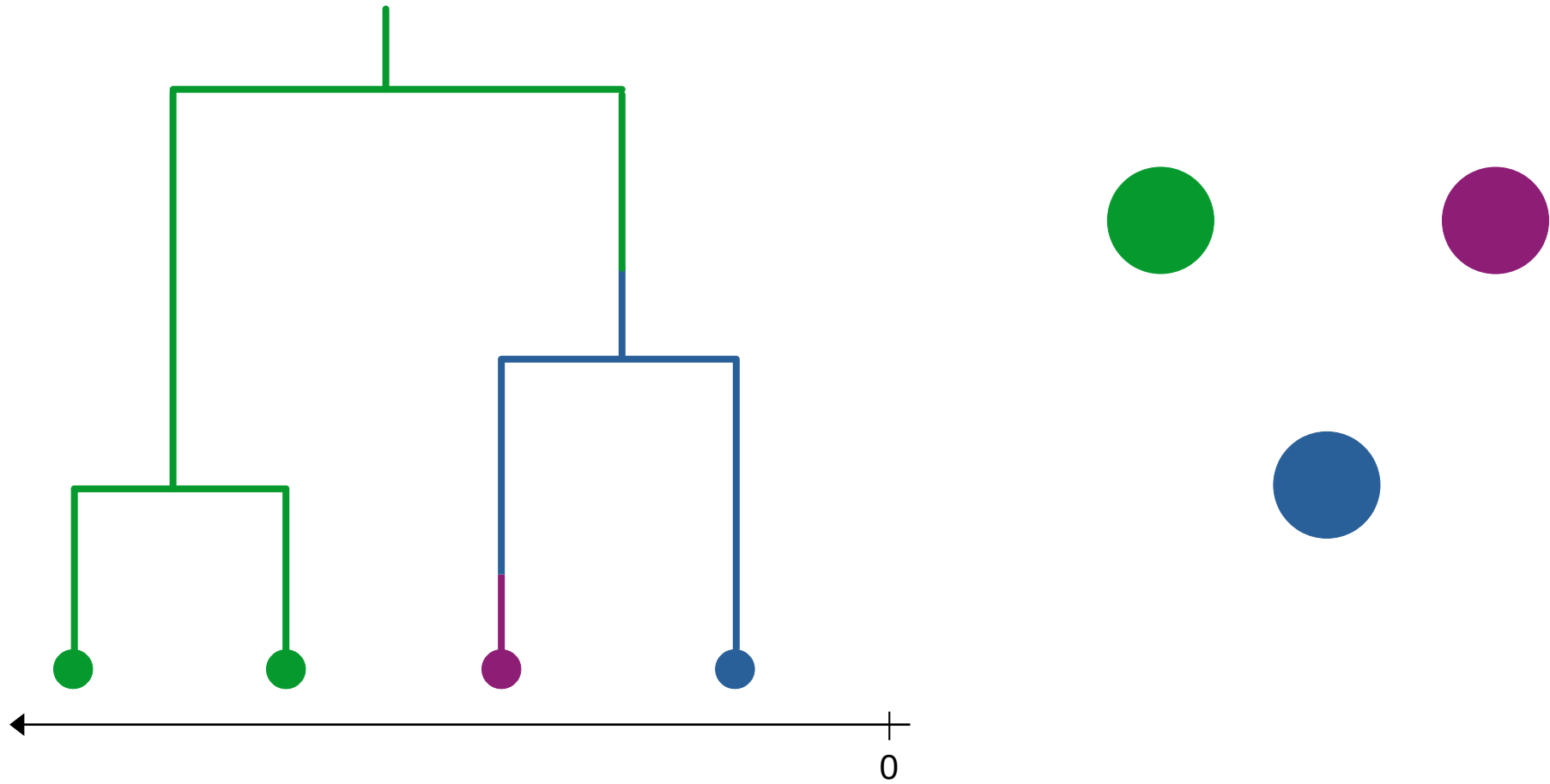




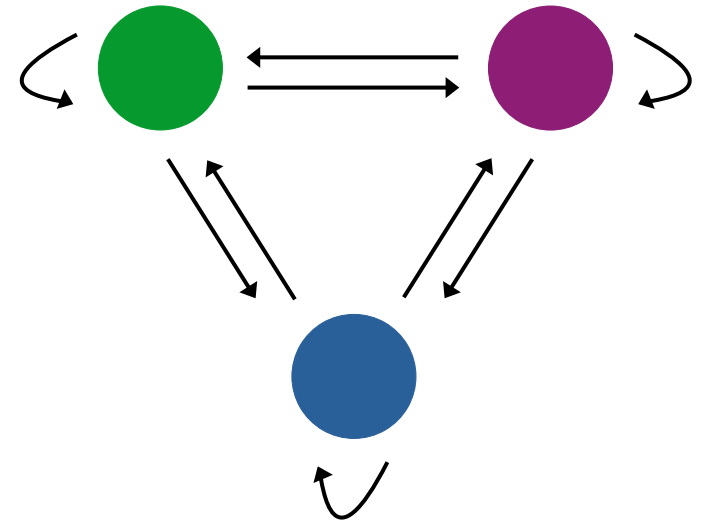
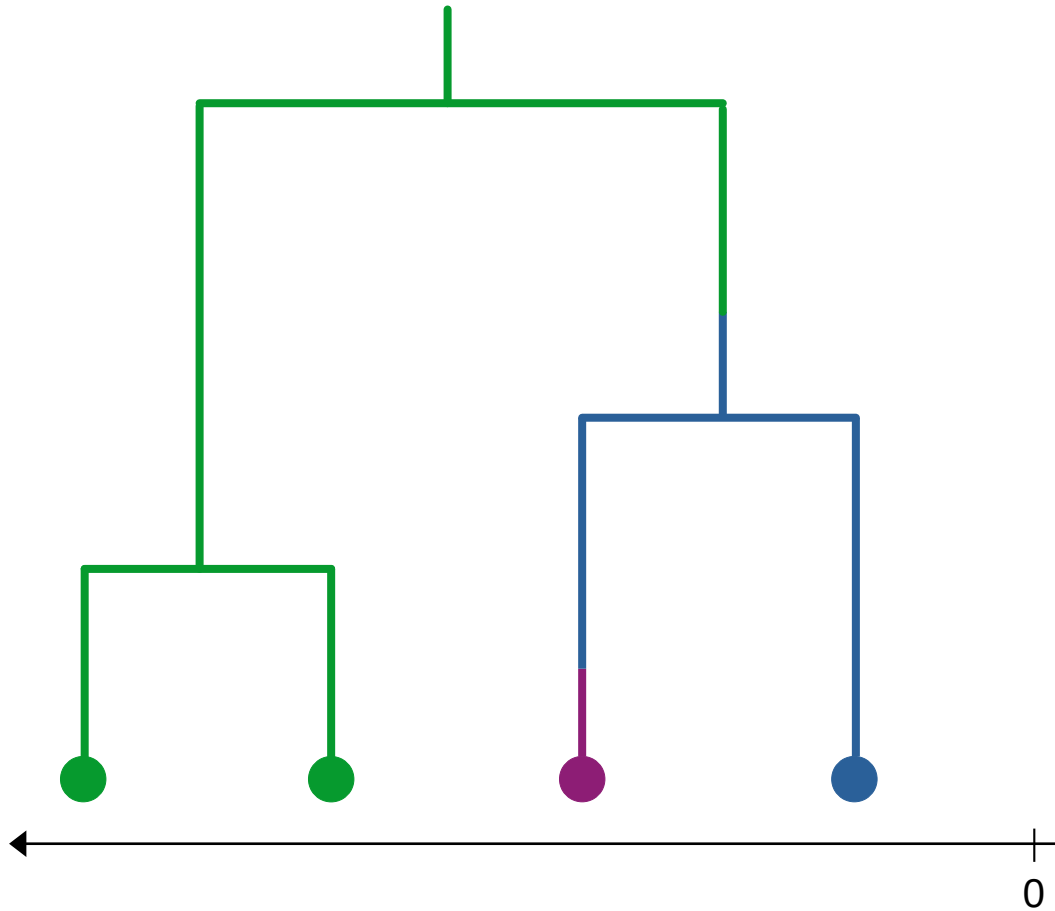
# Structured coalescent in action



# Structured coalescent in action



# Structured coalescent in action



# Example: source of a local outbreak

## French Foie Gras Makers Fear the Worst as Bird Flu Toll Rises

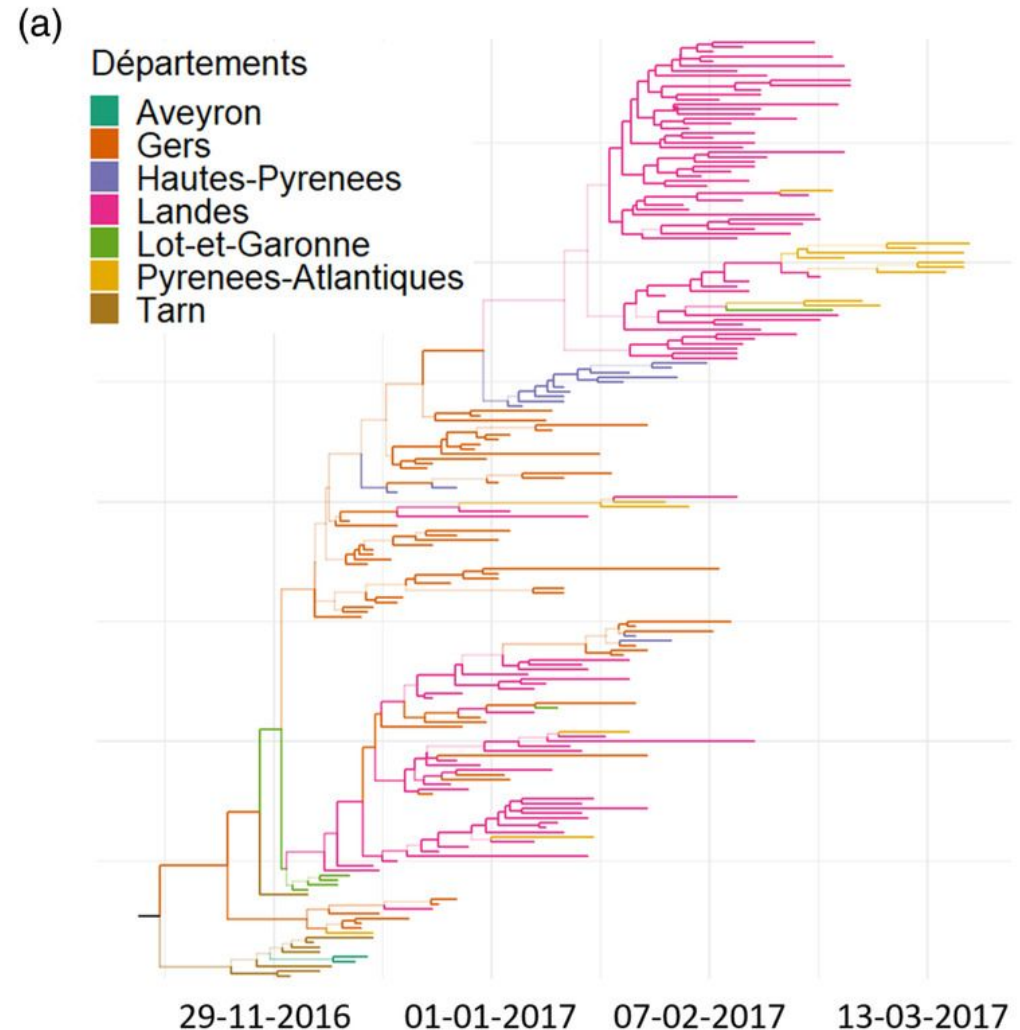
The bird flu epidemic in southwest France, home to most duck and foie gras producers, has led to more than three million poultry being killed.



Outbreak of Influenza A/H5N8 in the south of France 2016-2017

Tarn is inferred as the most likely origin of the outbreak

Chakraborty *et al.* **TBED 2022**



# In BEAST2



- MultiTypeTree (MTT): exact structured coalescent  
Vaughan *et al.* **Bioinformatics** 2014
- Structured COalescent Transmission Tree Inference (SCOTTI):  
approximate structured coalescent  
De Maio *et al.* **PLoS Genetics** 2015
  - assumes independence of lineages
  - assumes identical population sizes in subpopulations
- Marginal Approximation of the Structured Coalescent (MASCOT):  
approximate structured coalescent  
Müller *et al.* **Bioinformatics** 2018
  - assumes independence of lineages

# Key points to remember

- Structured models are designed to represent sub-populations within our dataset
- Structured birth-death models
  - Focus on differences in dynamics between sub-populations
  - Can be character-dependent or independent
  - Can be integrated with other BD models (skyline, FBD, etc)
- Structured coalescent models
  - Focus on the lack of interactions between sub-populations
  - Existing implementations are character-dependent

# BD vs coalescent: the revenge

- Clearer choice due to differences in assumptions and underlying process
- Both structured models are sensitive to sampling
  - BD relies on defined sampling process
  - Coalescent is sensitive to sampling biases
- Estimating migration rates
  - Disease outbreak: BD is more accurate
  - Endemic disease: both accurate, coalescent more precise

# Tutorial time



Character-dependent BD model (BDMM)

<https://taming-the-beast.org/tutorials/Structured-birth-death-model/>

Character-independent BD models (MSBD, ClaDS)

<https://taming-the-beast.org/tutorials/MSBD-tutorial/>

<https://taming-the-beast.org/tutorials/ClaDS-tutorial/>

Structured coalescent (MASCOT)

<https://taming-the-beast.org/tutorials/Mascot-Tutorial/>