

# Phylogenetics

Introduction to  
phyldynamics models  
RL-V3 MPP

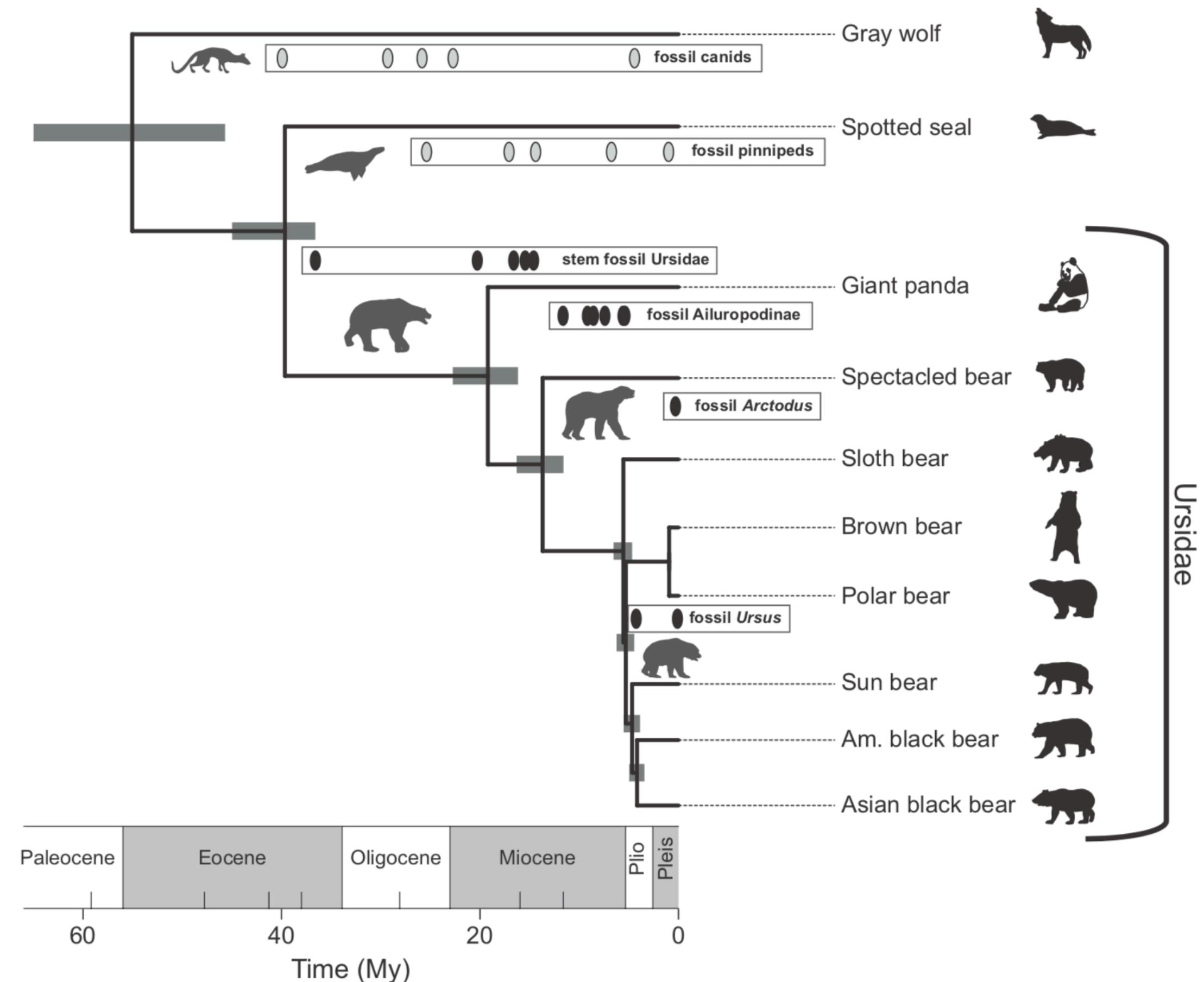
Rachel Warnock

16.04.2025



# Objectives

- Recap: tripartite framework
- The fossilised birth-death process
- Total-evidence dating
- Phylodynamics



# Bayesian divergence time estimation

Recap

# We use a Bayesian framework

$$P(\text{ model } | \text{ data }) = \frac{P(\text{ data } | \text{ model }) P(\text{ model })}{P(\text{ data })}$$

likelihood

priors

posterior

marginal probability of the data

# Bayesian divergence time estimation

## The data

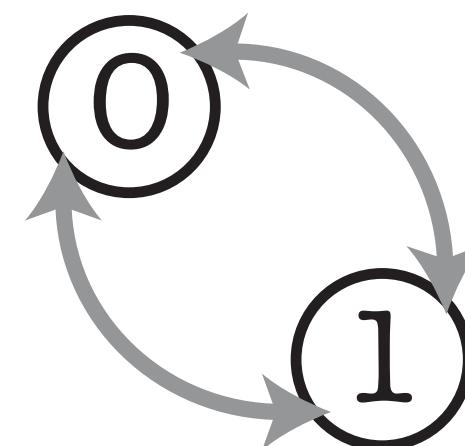
and / or  
0101... ATTG...  
1101... TTGC...  
0100... ATTC...



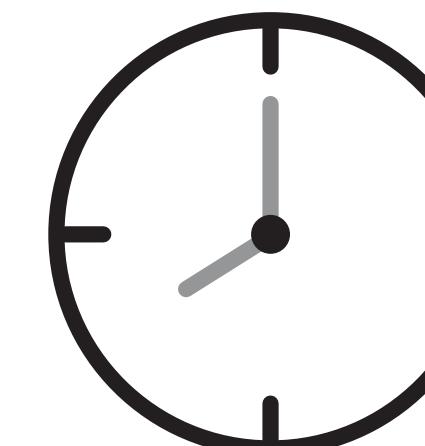
phylogenetics  
characters

sample  
ages

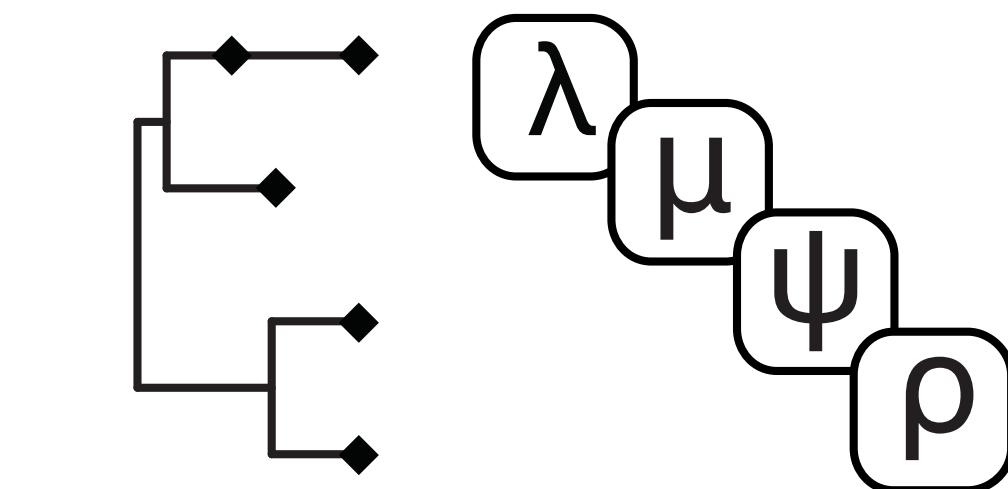
## 3 model components



substitution  
model



clock  
model



tree and tree  
model

# Bayesian divergence time estimation

posterior

$$P(E \mid \lambda, \mu, \psi, p, O, t \mid 0101\dots, 1101\dots, 0100\dots, \text{snail}) =$$

likelihood

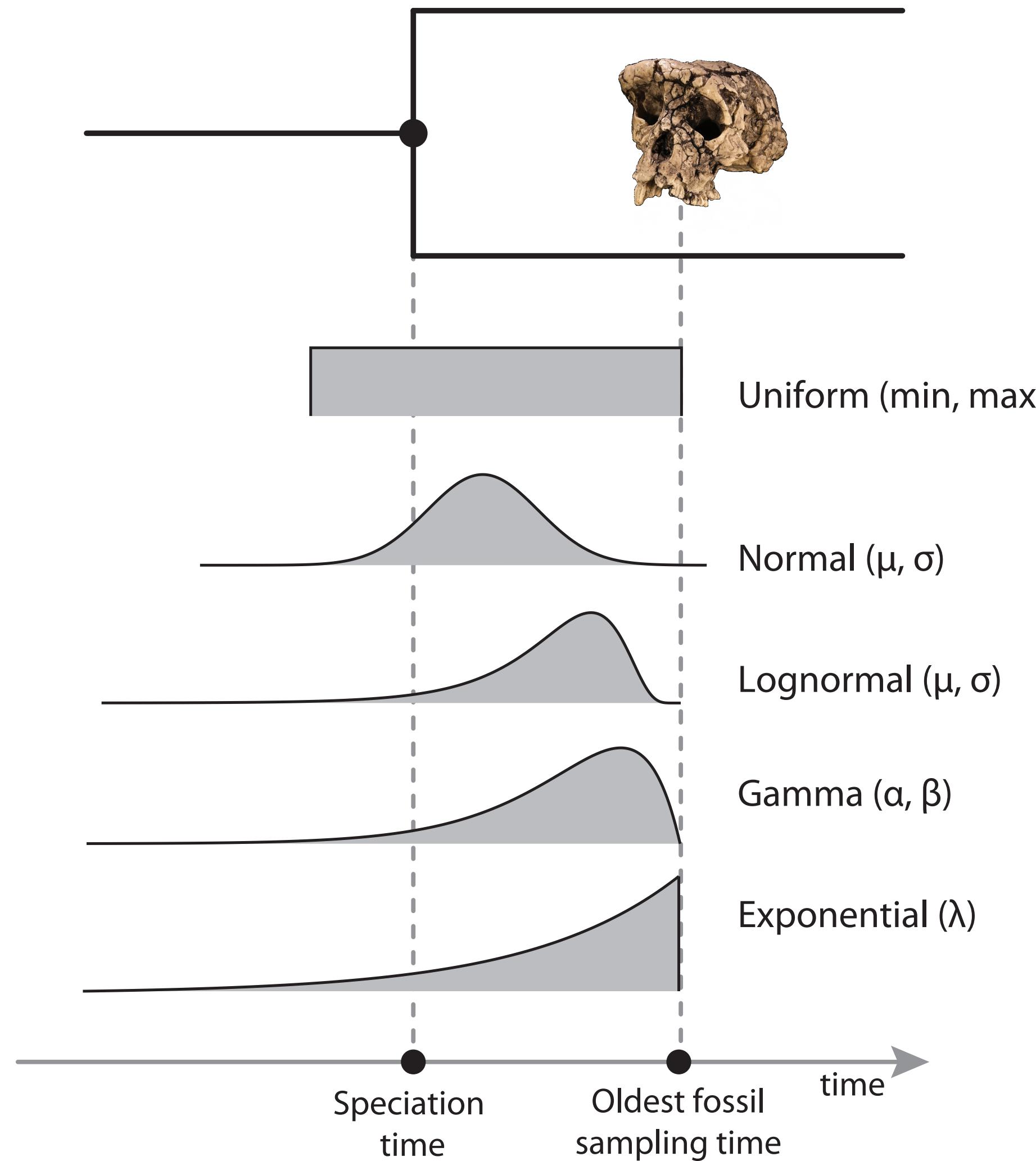
probability of the  
time tree

priors

$$P(0101\dots \mid E) P(E \mid \lambda, \mu, \psi, p, O, t) P(\lambda, \mu, \psi, p) P(O) P(t)$$
$$P(0101\dots \mid \text{snail})$$

marginal pr of the data

# Recap: Node dating



We can use a **calibration density** to constrain internal node ages

We typically use a **birth-death process** model to describe the tree generating process

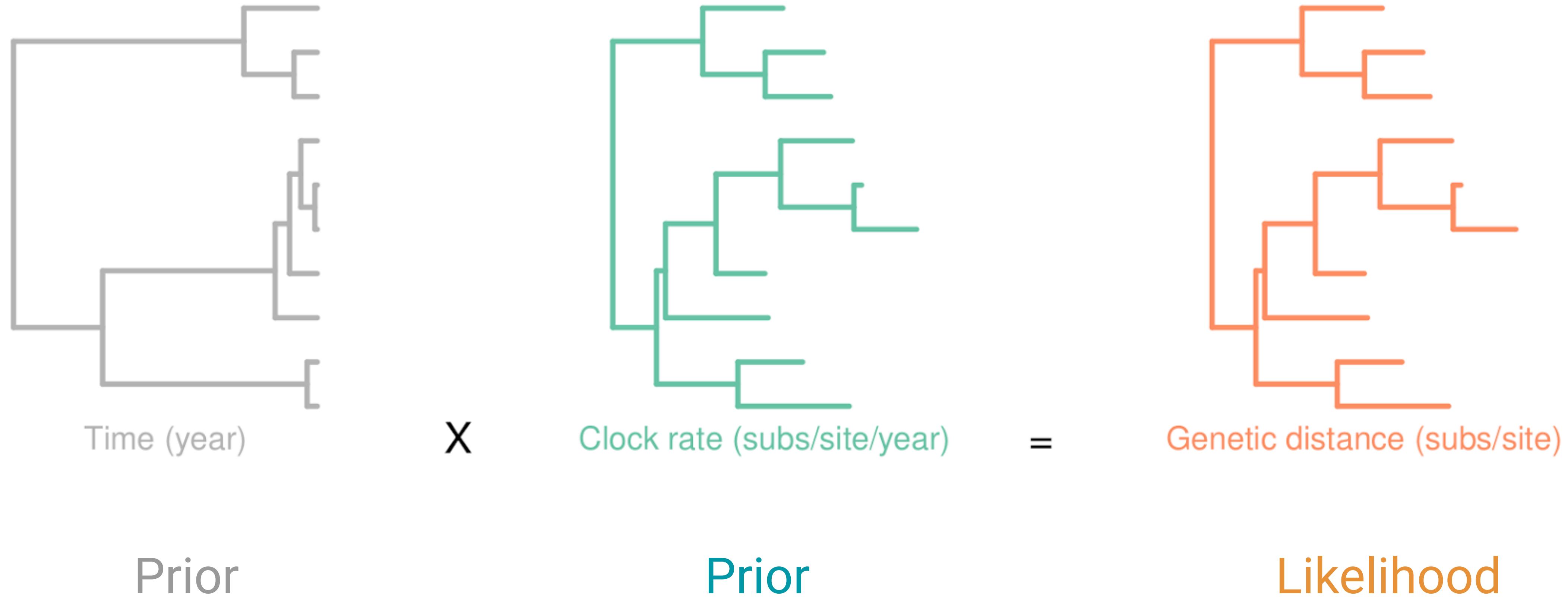
Adapted from Heath (2012). Sys Bio

# Issues with node dating

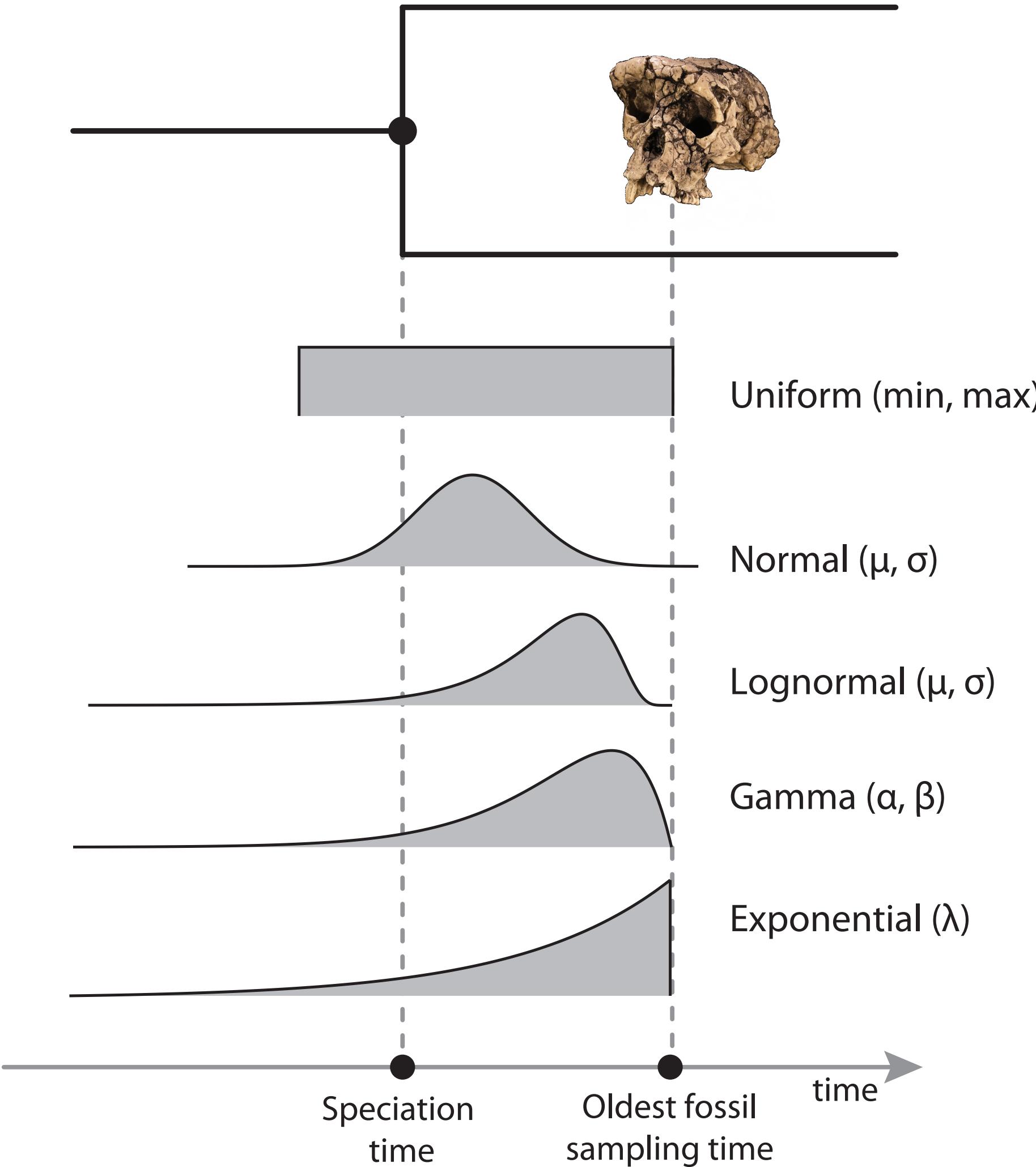
1. Rate and time is non-identifiable
2. Calibrations are hard to define objectively
3. The “effective priors” do not match the user specified priors
4. We’re potentially excluding a lot of information

# The priors will always influence our results

Rate and time are not identifiable - not a conventional Bayesian problem



# Minimum and maximum constraints



**Hard** minimum bounds are based on first appearances

**Soft** maximum bounds are based on more tenuous evidence – typically the 97.5% limit of the calibration density

The upper bounds often influence the results

# Maximum constraints

## Best Practices for Justifying Fossil Calibrations

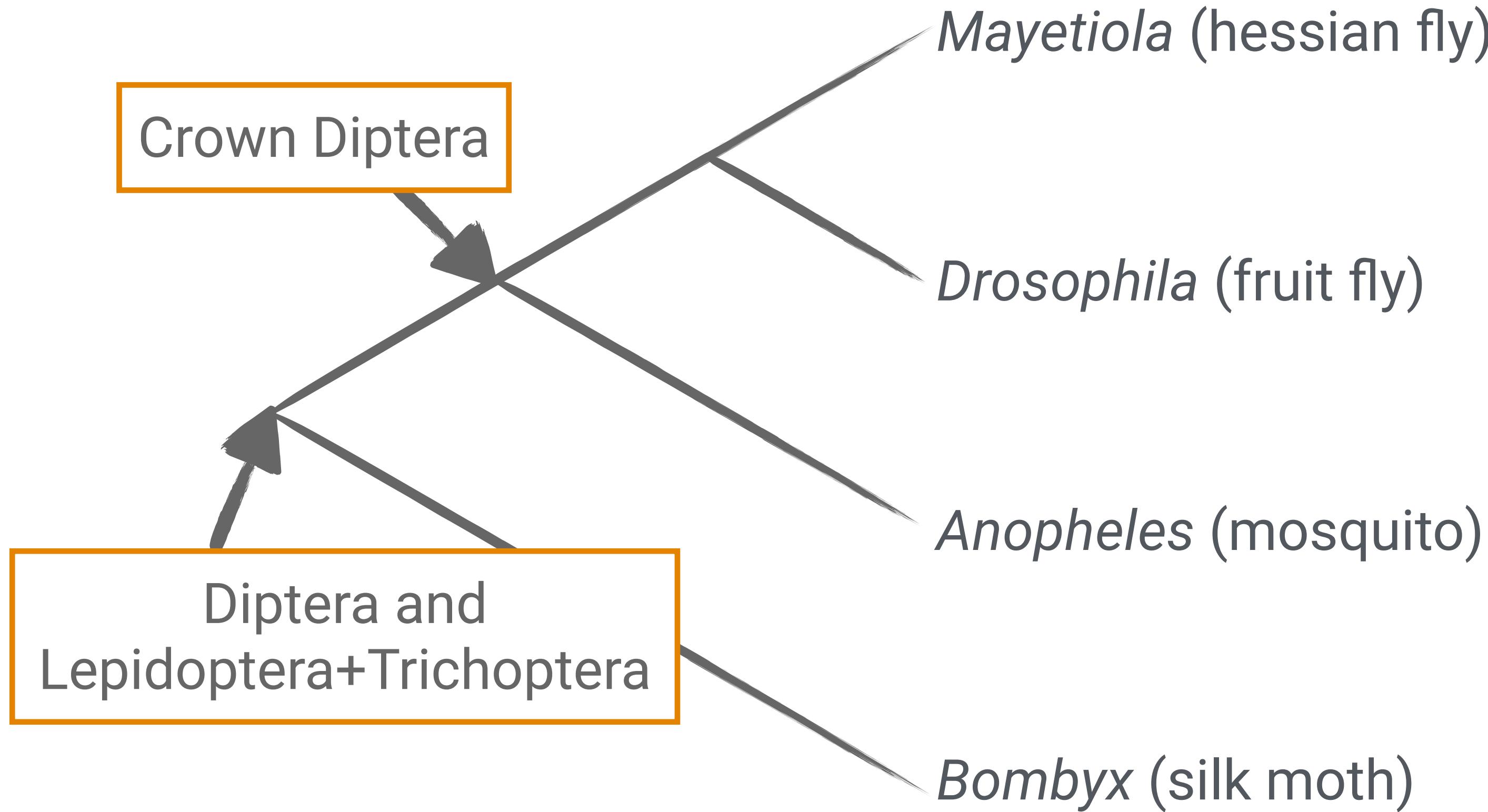
JAMES F. PARHAM<sup>1,2,\*</sup>, PHILIP C. J. DONOGHUE<sup>3</sup>, CHRISTOPHER J. BELL<sup>4</sup>, TYLER D. CALWAY<sup>5</sup>,  
JASON J. HEAD<sup>6</sup>, PATRICIA A. HOLROYD<sup>7</sup>, JUN G. INOUE<sup>8</sup>, RANDALL B. IRMIS<sup>9</sup>,  
WALTER G. JOYCE<sup>10</sup>, DANIEL T. KSEPKA<sup>11,12</sup>, JOSÉ S. L. PATANÉ<sup>13</sup>, NATHAN D. SMITH<sup>14,15</sup>,  
JAMES E. TARVER<sup>3,16</sup>, MARCEL VAN TUINEN<sup>17</sup>, ZIHENG YANG<sup>18</sup>, KENNETH D. ANGIELCZYK<sup>15</sup>,  
JENNY M. GREENWOOD<sup>3</sup>, CHRISTY A. HIPSLEY<sup>19,20</sup>, LOUIS JACOBS<sup>21</sup>, PETER J. MAKOVICKY<sup>15</sup>,  
JOHANNES MÜLLER<sup>19</sup>, KRISTER T. SMITH<sup>22</sup>, JESSICA M. THEODOR<sup>23</sup>, RACHEL C. M. WARNOCK<sup>3</sup>,  
AND MICHAEL J. BENTON<sup>3</sup>

Goal - to make calibration choices transparent and explicit

*"The [soft] maximum constraint is established as older than all the oldest possible records, extending back to encompass a time when the ecologic, biogeographic, geologic, and taphonomic conditions for the existence of the lineage are met, but no records are known."*

# Min & max constraints

Example from insects



## Minimum (all nodes)

- 238.5 Ma
- Triassic Grès-a-Voltzia Frm, France
- Earliest (non-controversial) evidence for all 4 lineages

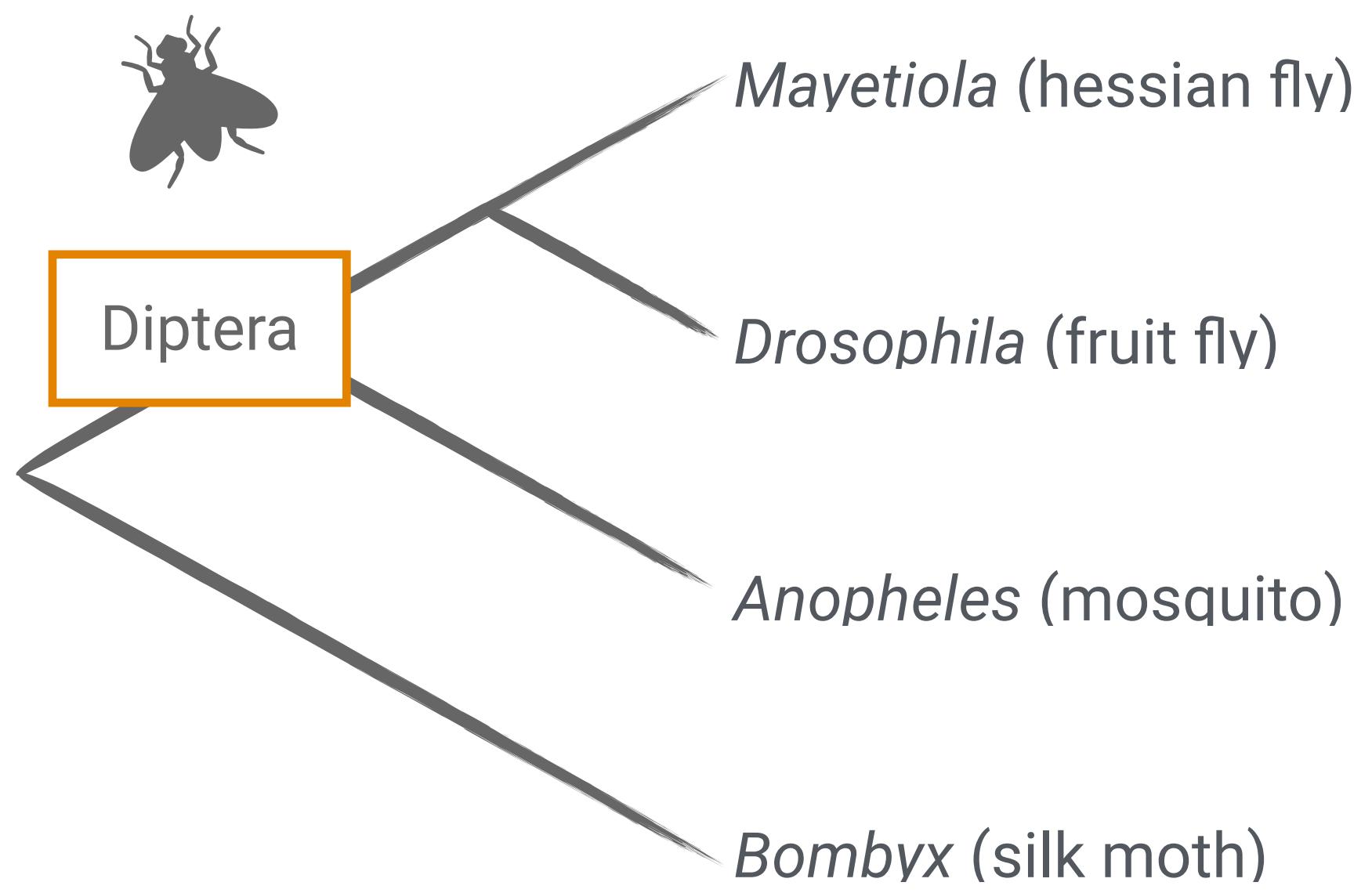
## Maximum (all nodes)

- 295.4 Ma
- Boskovice Furrow, Moravia, Czechia
- Huge diversity of insects described from here - no members of even total group Diptera from here or younger deposits

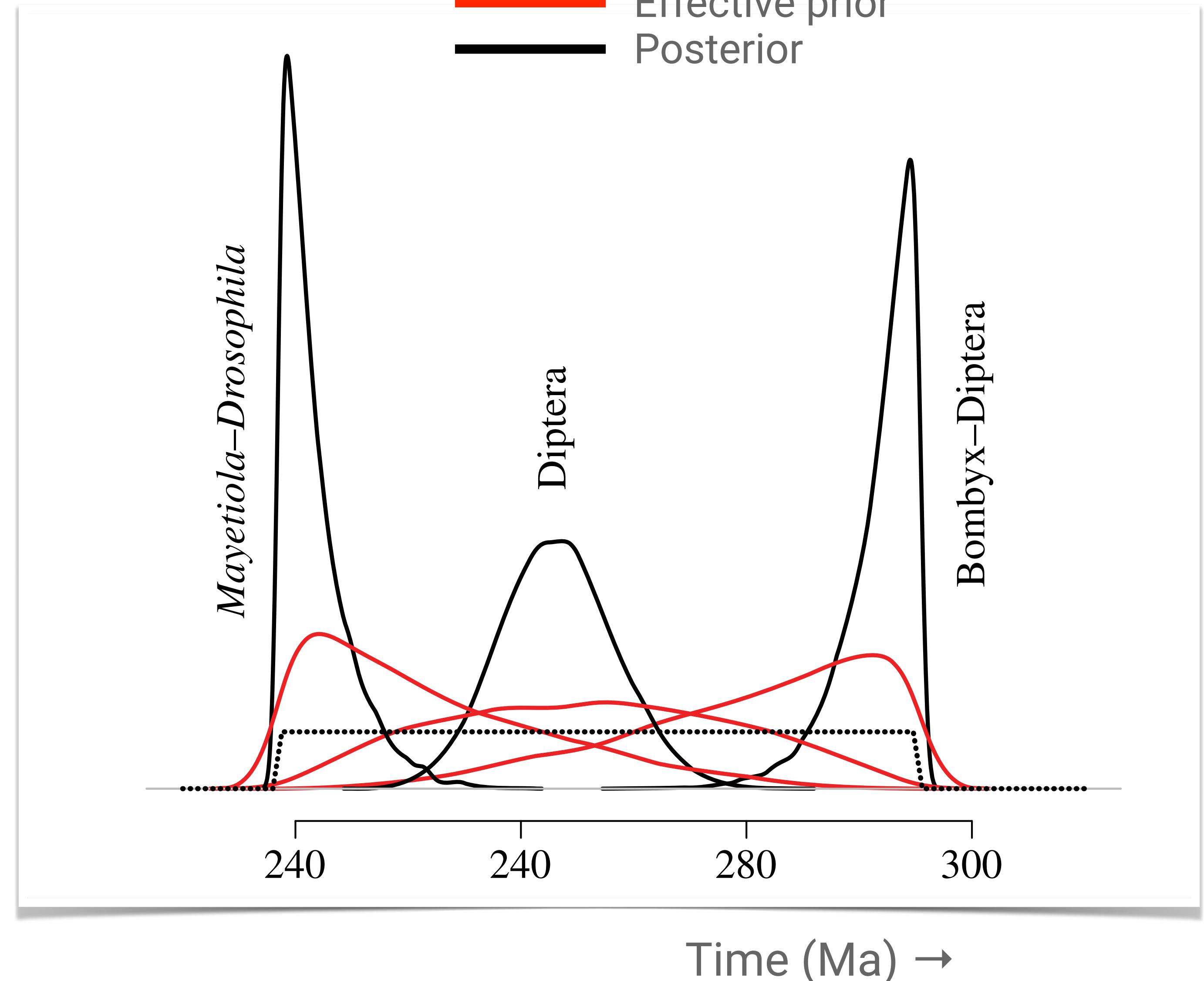
# The user “specified priors” will not (always) match the “effective priors” used during analysis



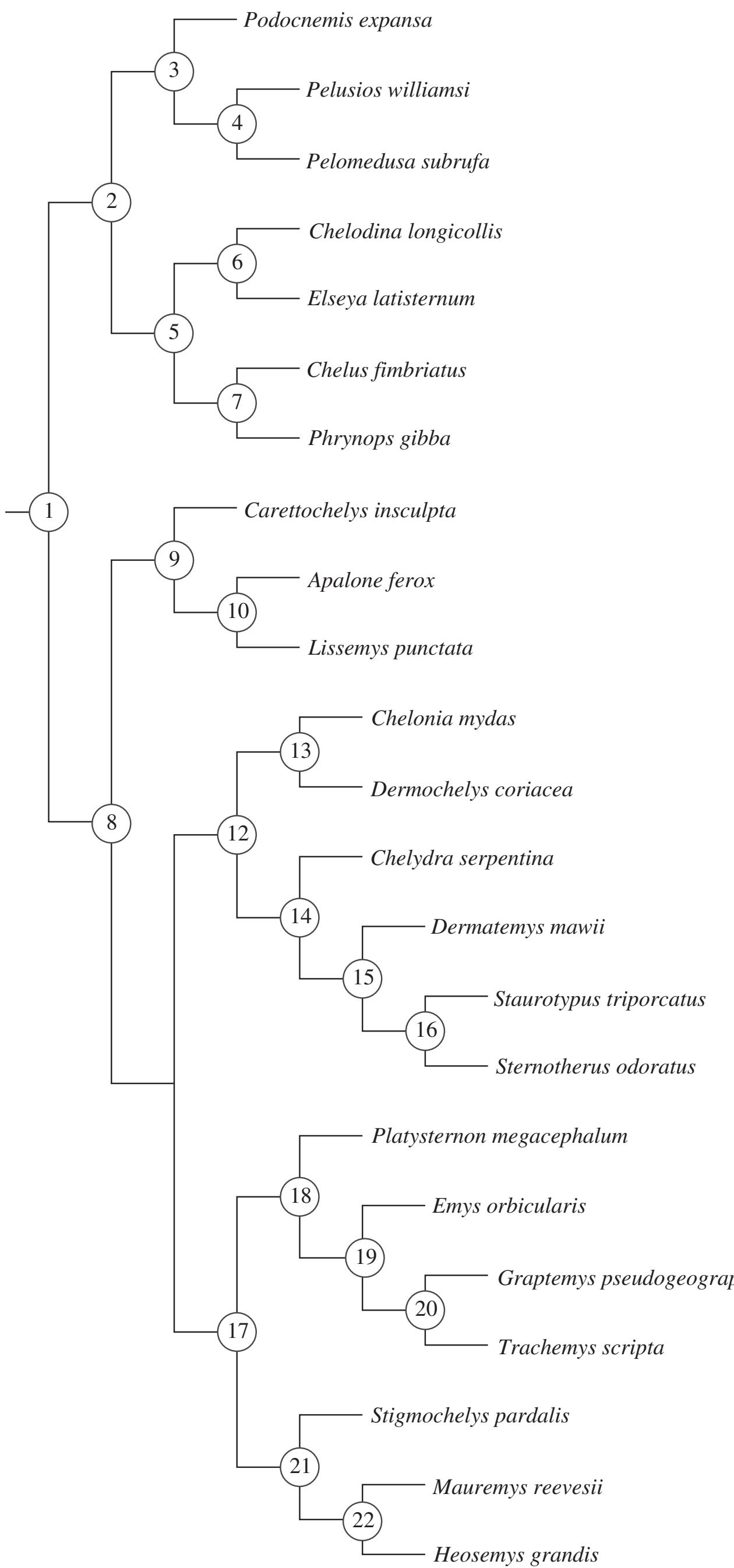
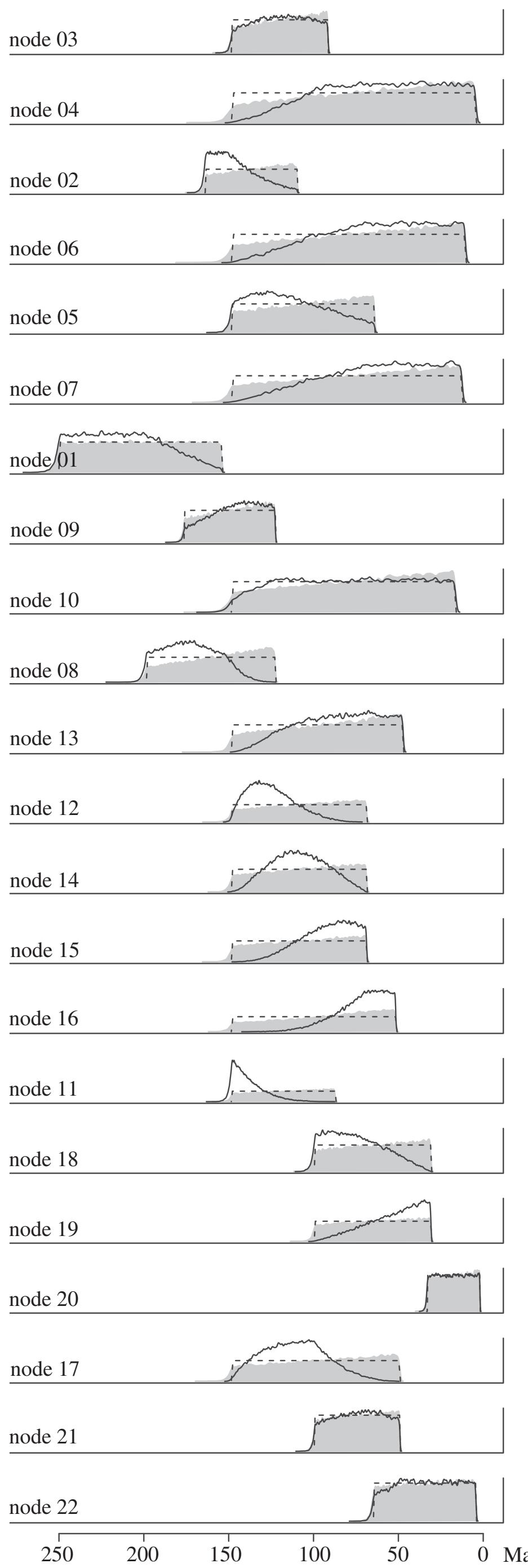
We can run our analysis “under the prior” (we ignore the sequence data)  
to see how the software actually constructs the prior density  
– accounts for the interaction between nodes



Posteriors look suspiciously like  
the effective priors....



3.



## An example from turtles



Dashed lines = the specified priors

Grey shaded area = effective priors  
for one calibration, analysed alone

Black shaded area = effective  
priors for all calibrations together

The model doesn't describe the process that generated the fossil sampling times, meaning the model is **statistically incoherent**

The calibration priors are difficult to specify objectively and can have a massive impact on the divergence times. They can also interact with each other and / or the birth-death process prior in unintuitive ways

Some references on issues with specified vs effective priors

Yang and Rannala. [2006](#). MBE

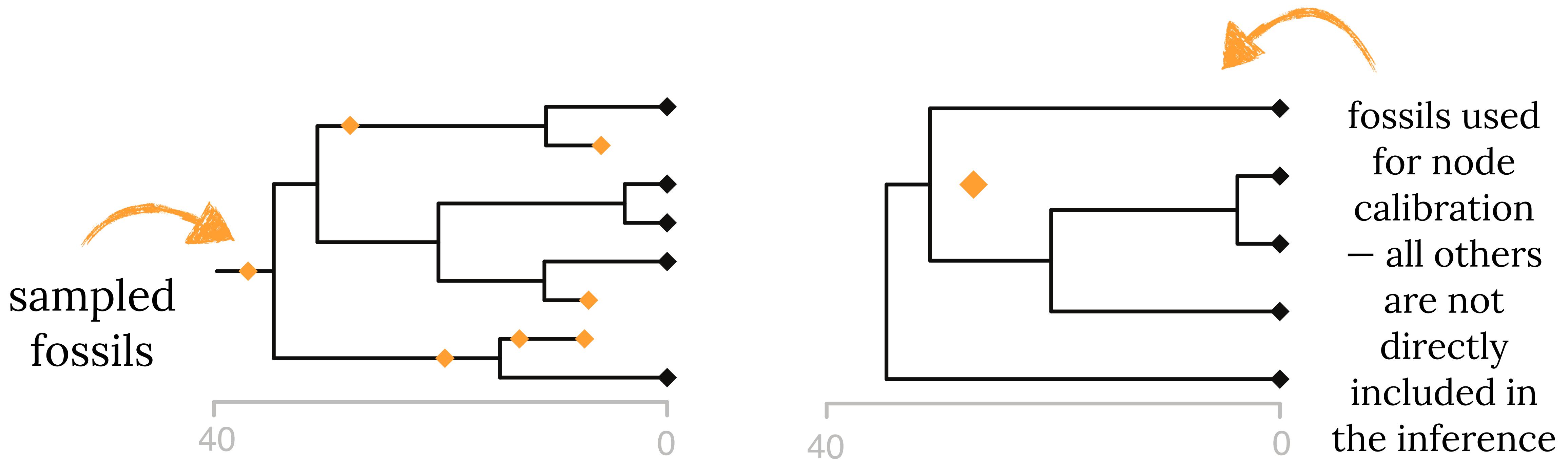
Heled and Drummond. [2012](#). Sys Bio

Warnock et al. [2012, 2015](#)

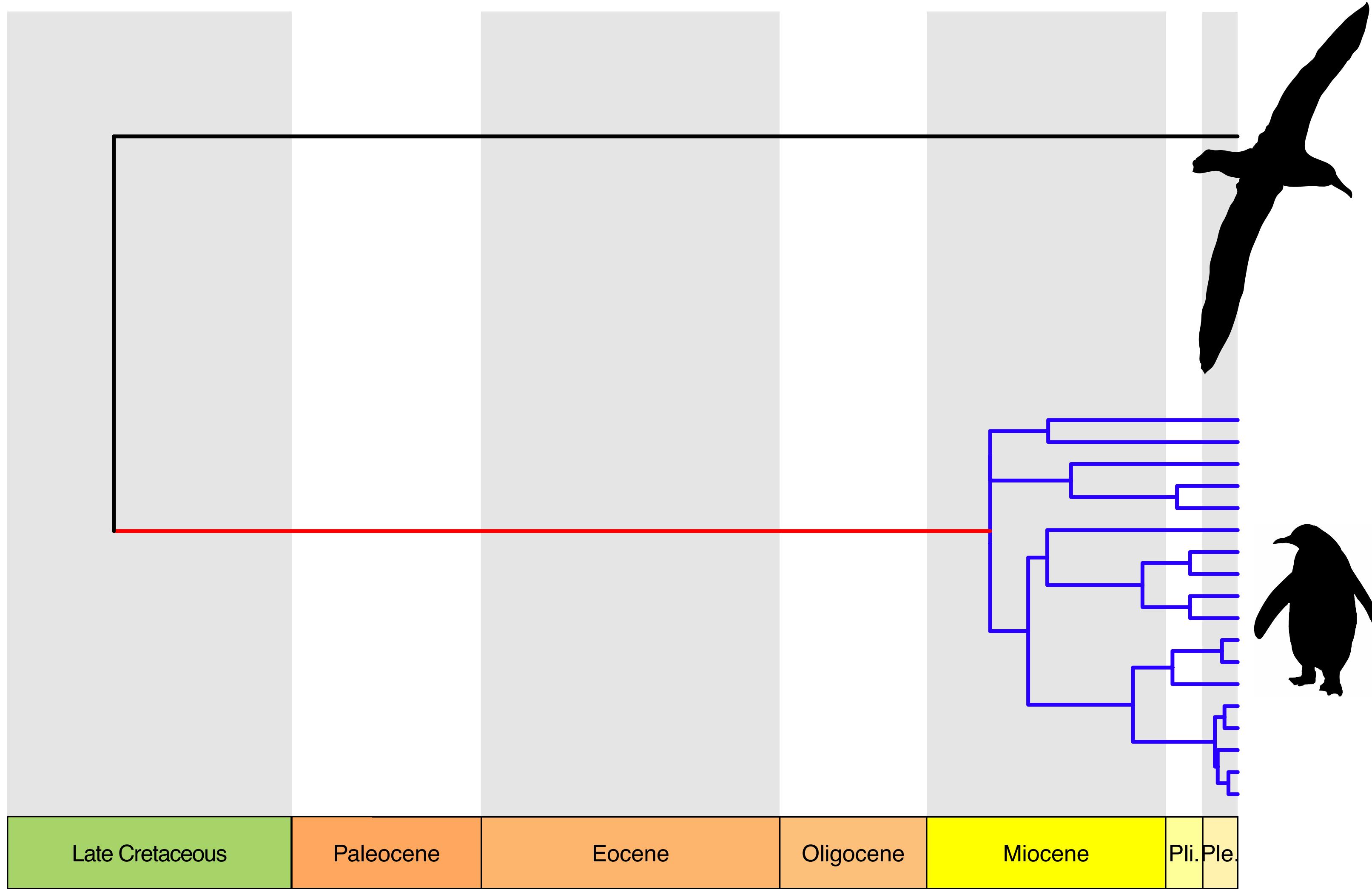
# Node dating: potential issues

There are many!

A lot of information is excluded, since typically we assign one fossil per calibration node

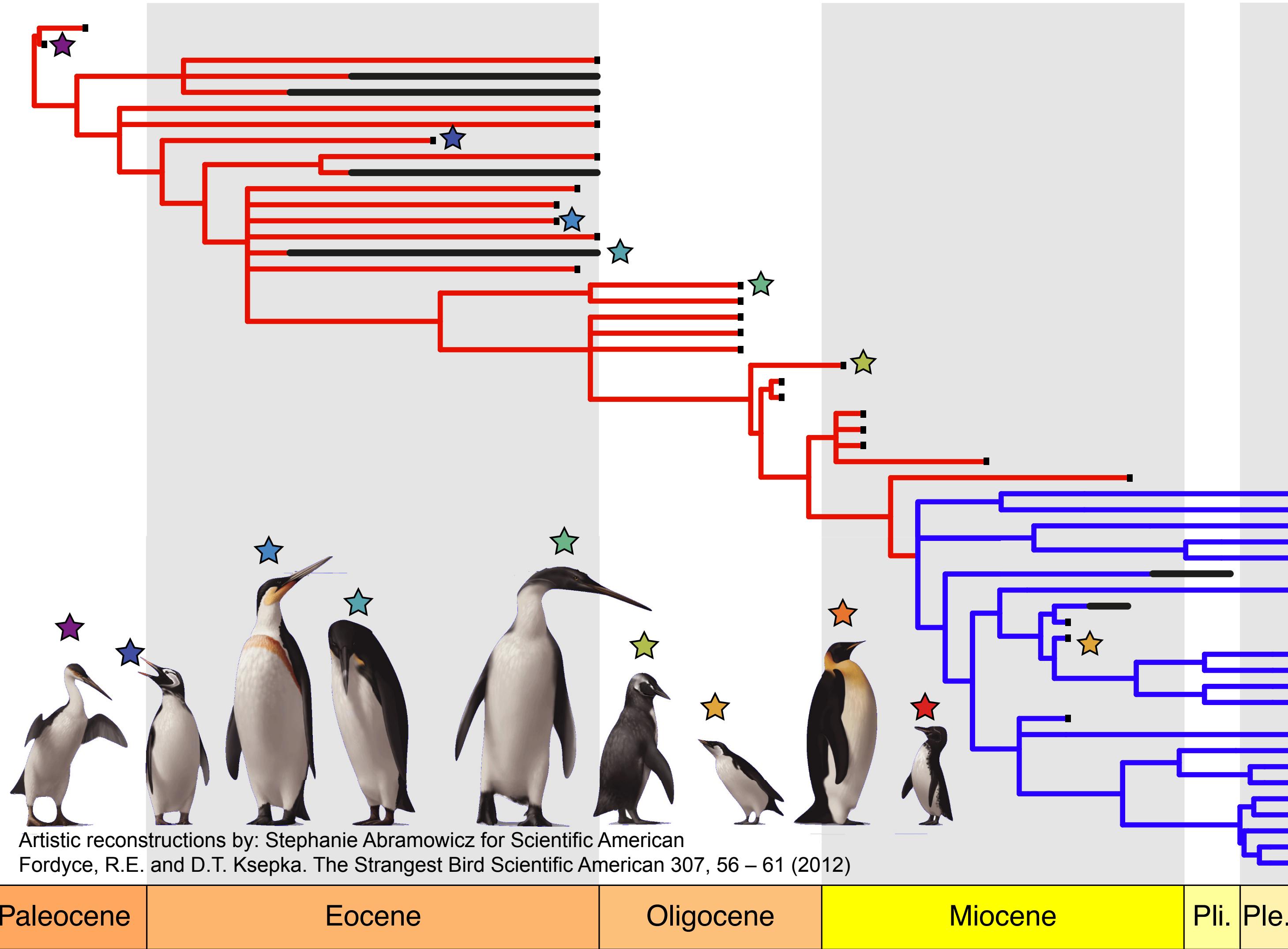


# Example: living penguins



Nearest living relative is the group containing falcons - separated by ~60 Ma

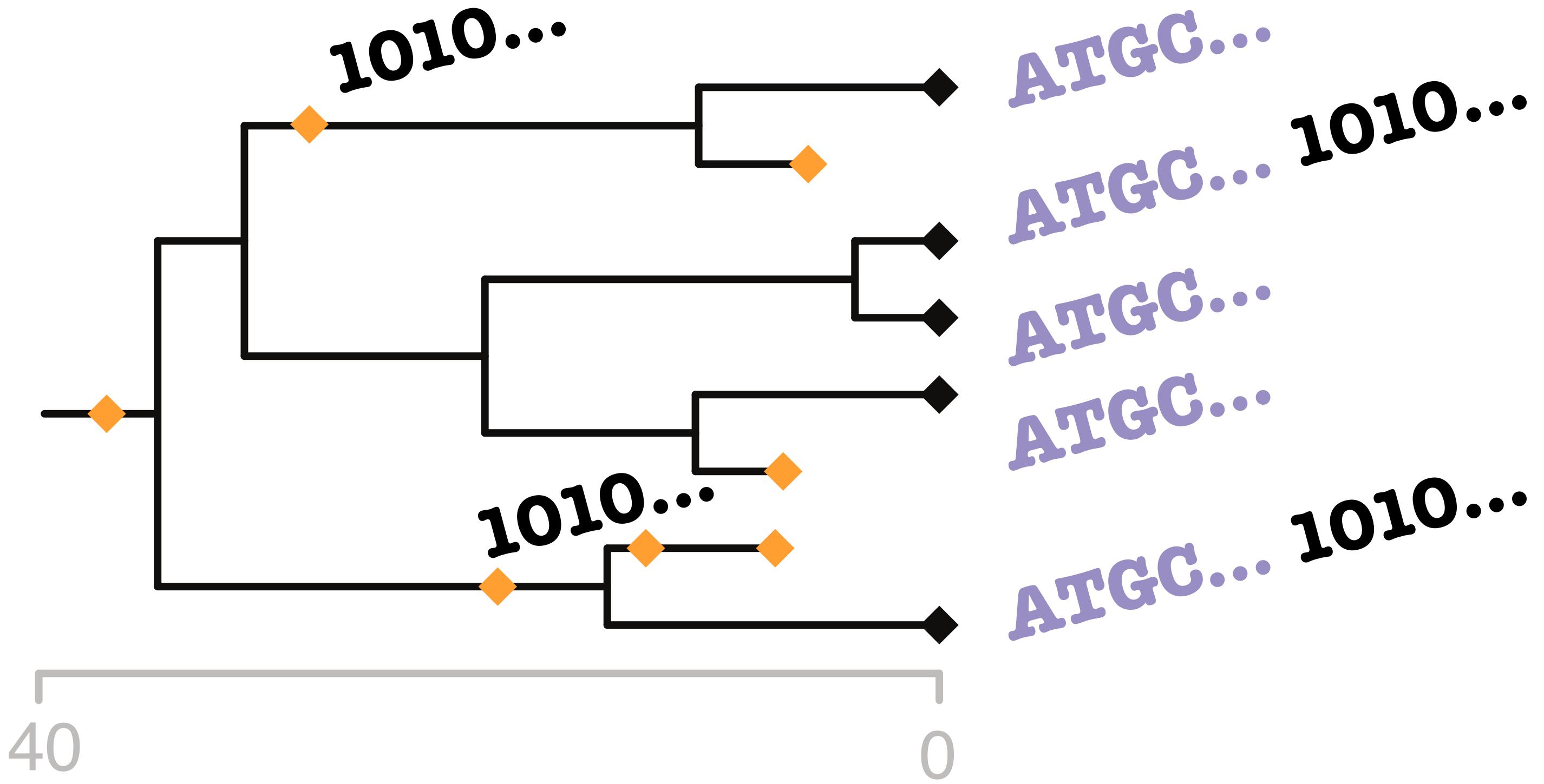
# Example: living penguins



But penguins  
have a rich  
fossil record!

# Total-evidence dating

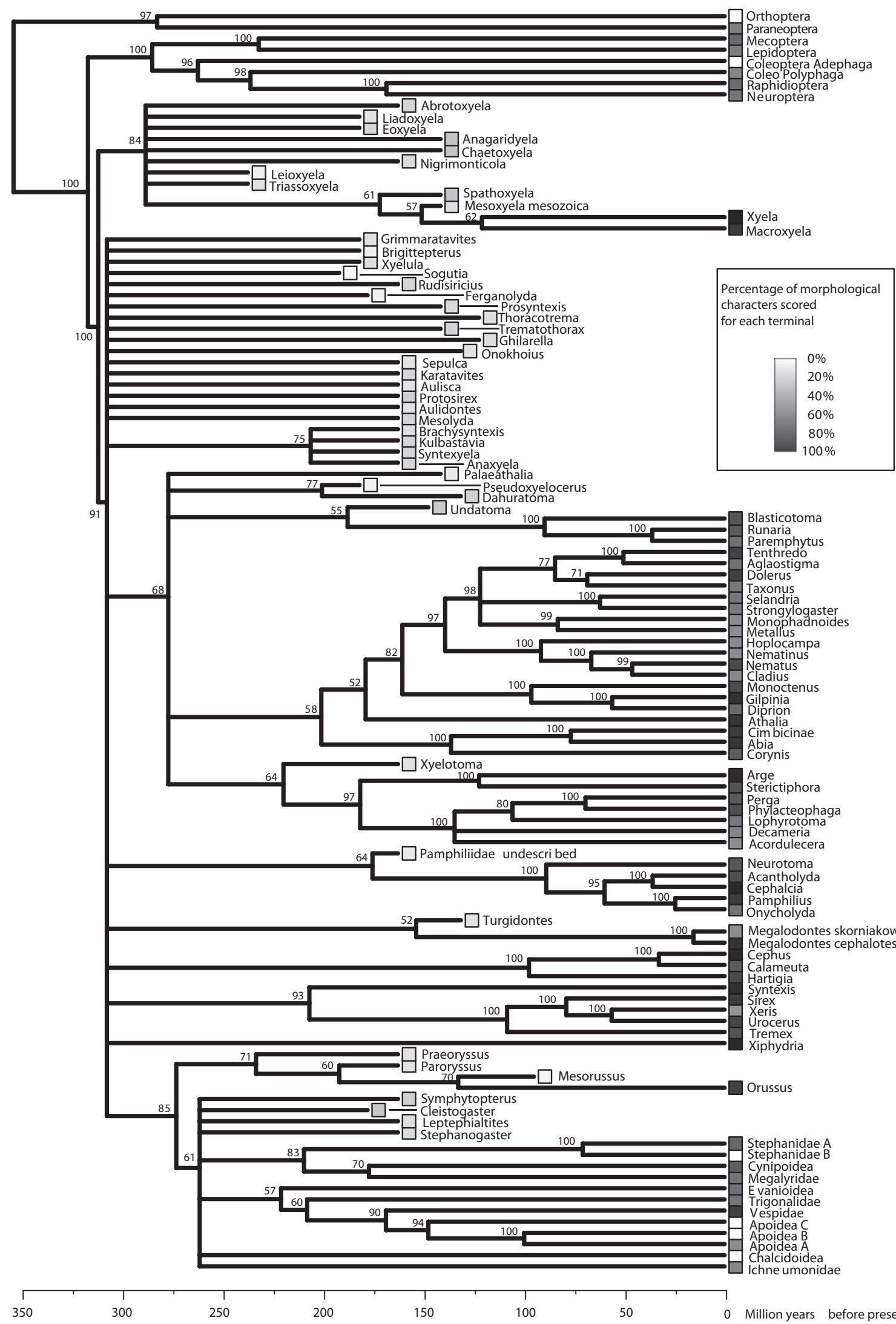
# Tip-dating or “total-evidence” dating



We have DNA for living species. We have morphology for living *and* fossil species

Fossils can be positioned on the basis of morphology  
→ accounts for uncertainty in fossil placement

# The uniform tree prior

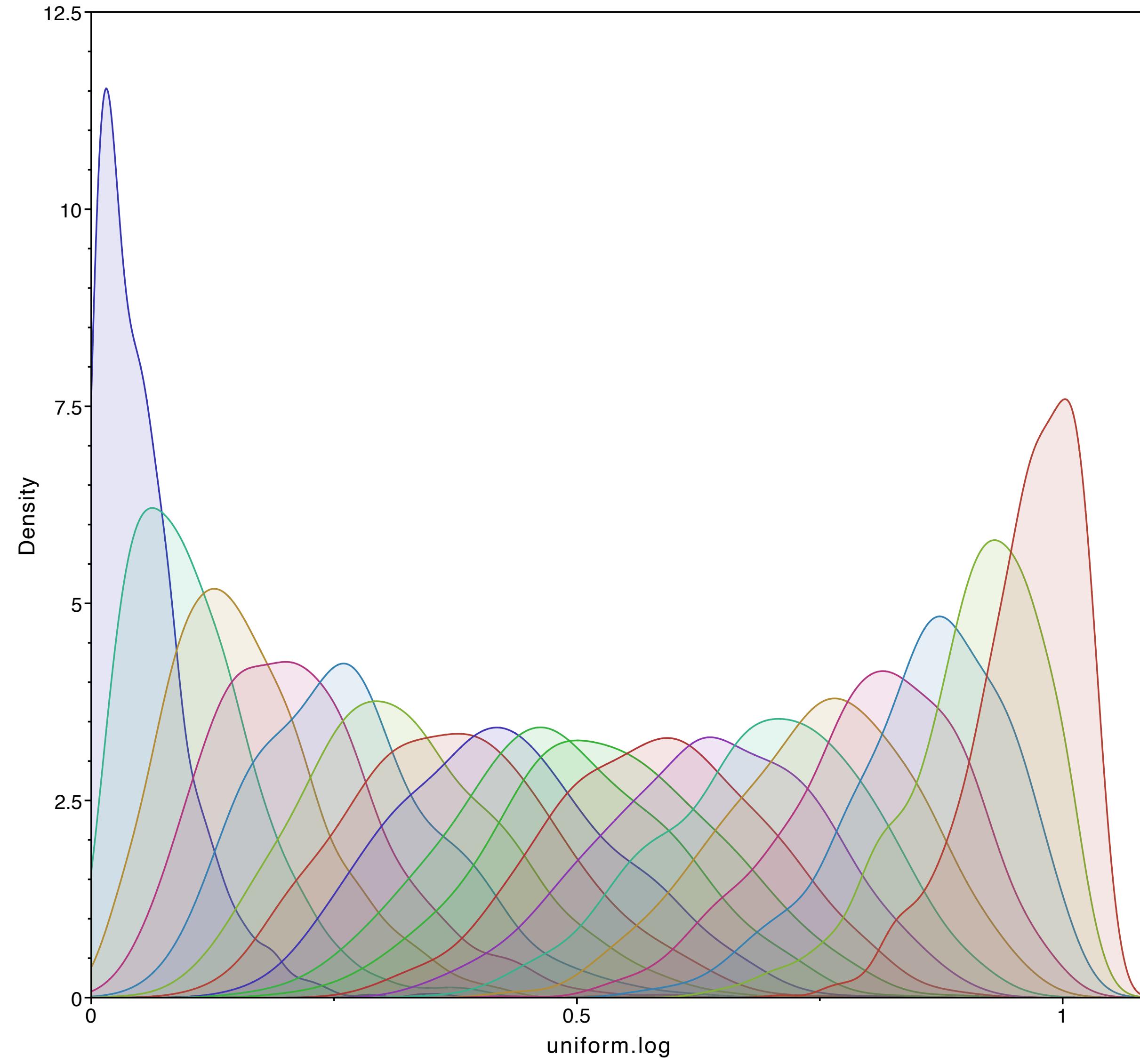


Dated tree of Hymenoptera



The uniform tree prior assumes all trees and branch lengths are equally likely within the bounds of the fossil ages (+ a max upper bound)

It does not explicitly account for the fossil sampling process



A uniform tree prior implies time till the next split is independent of how many lineages there are present

This is in contrast to birth-death processes, where more lineages mean a higher chance of observing a split in one of these lineages

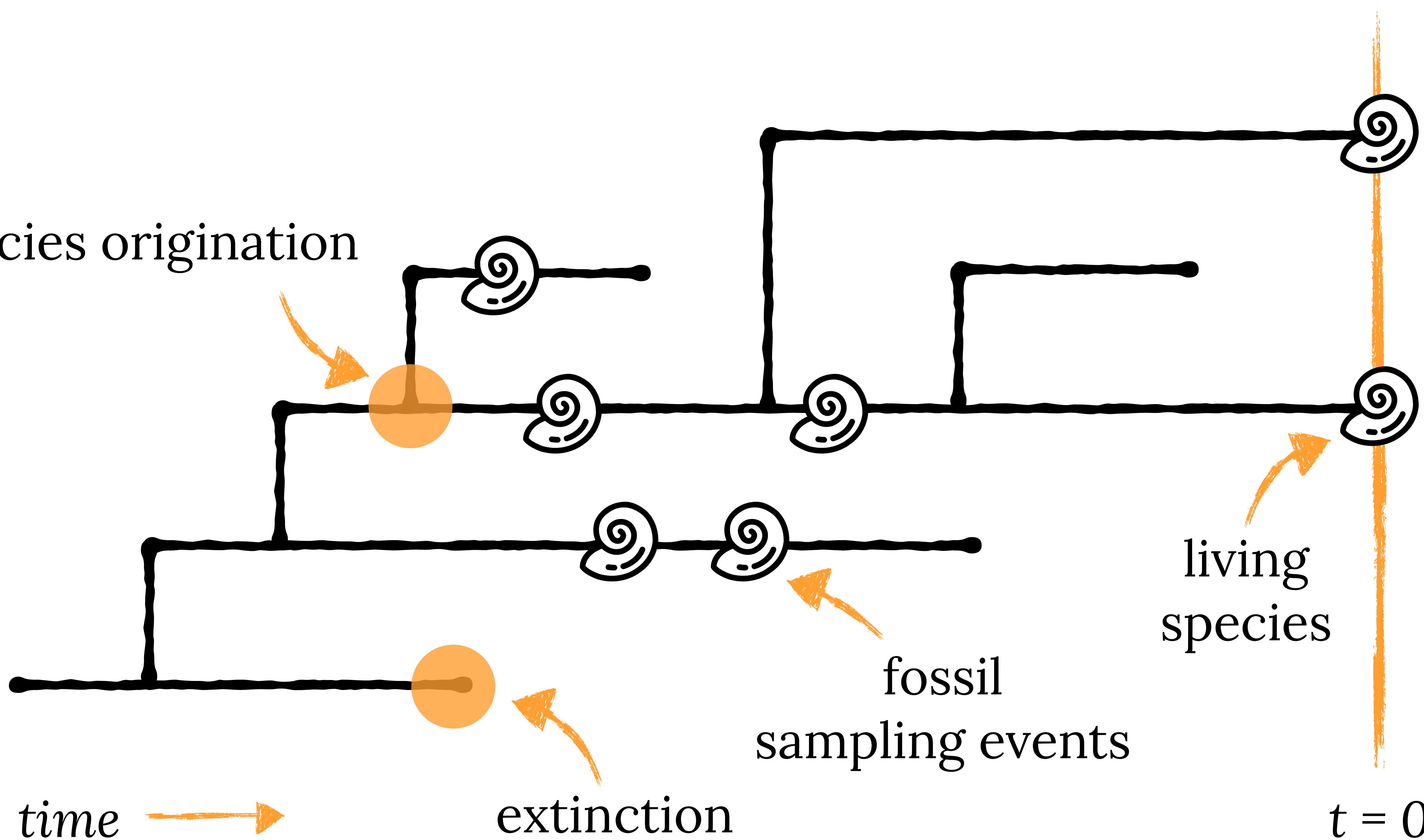
[Uniform tree priors - why not use them?](#)

Remco Bouckaert

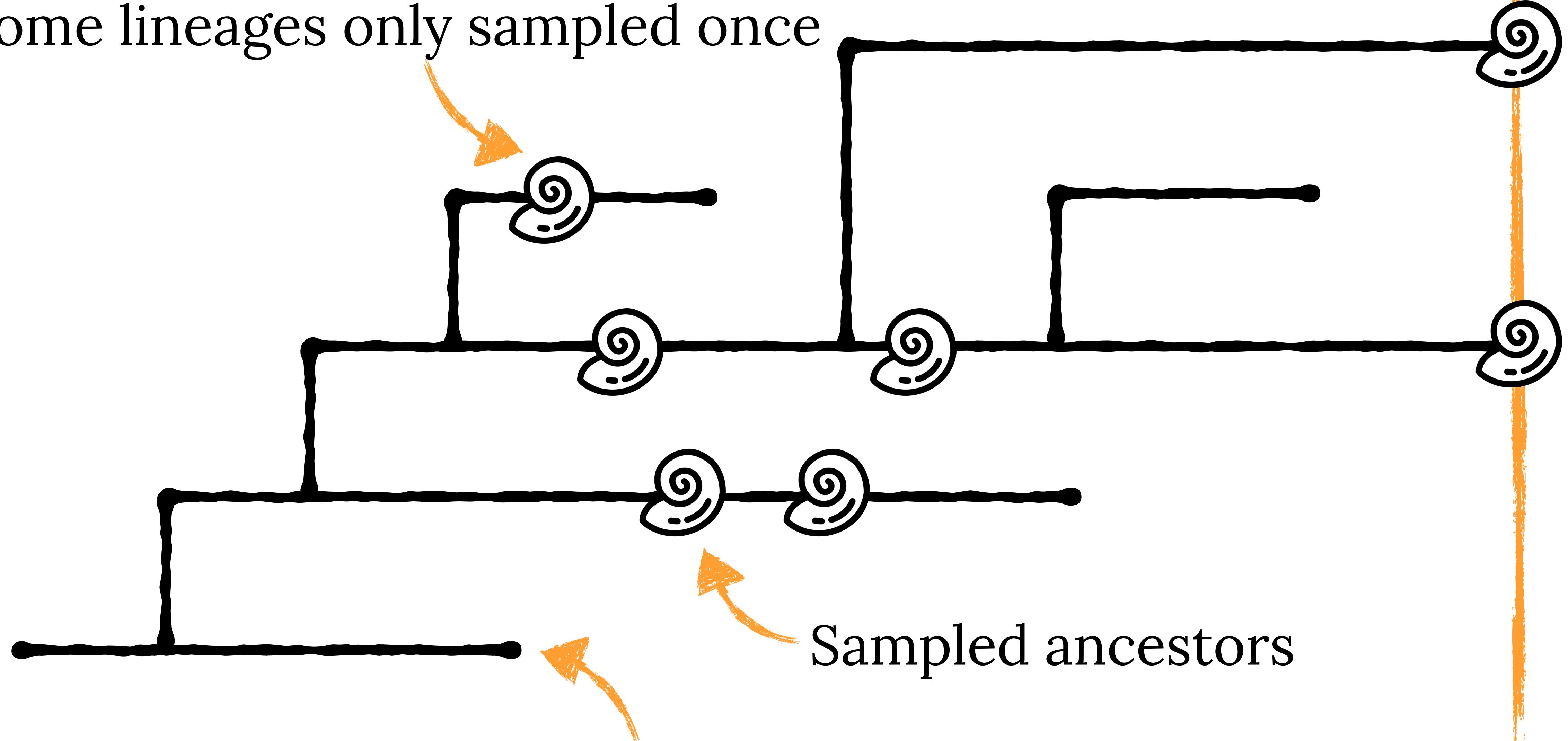
What does a generating prior for the fossil record look like?

# The fossilised birth-death process

species origination



Some lineages only sampled once

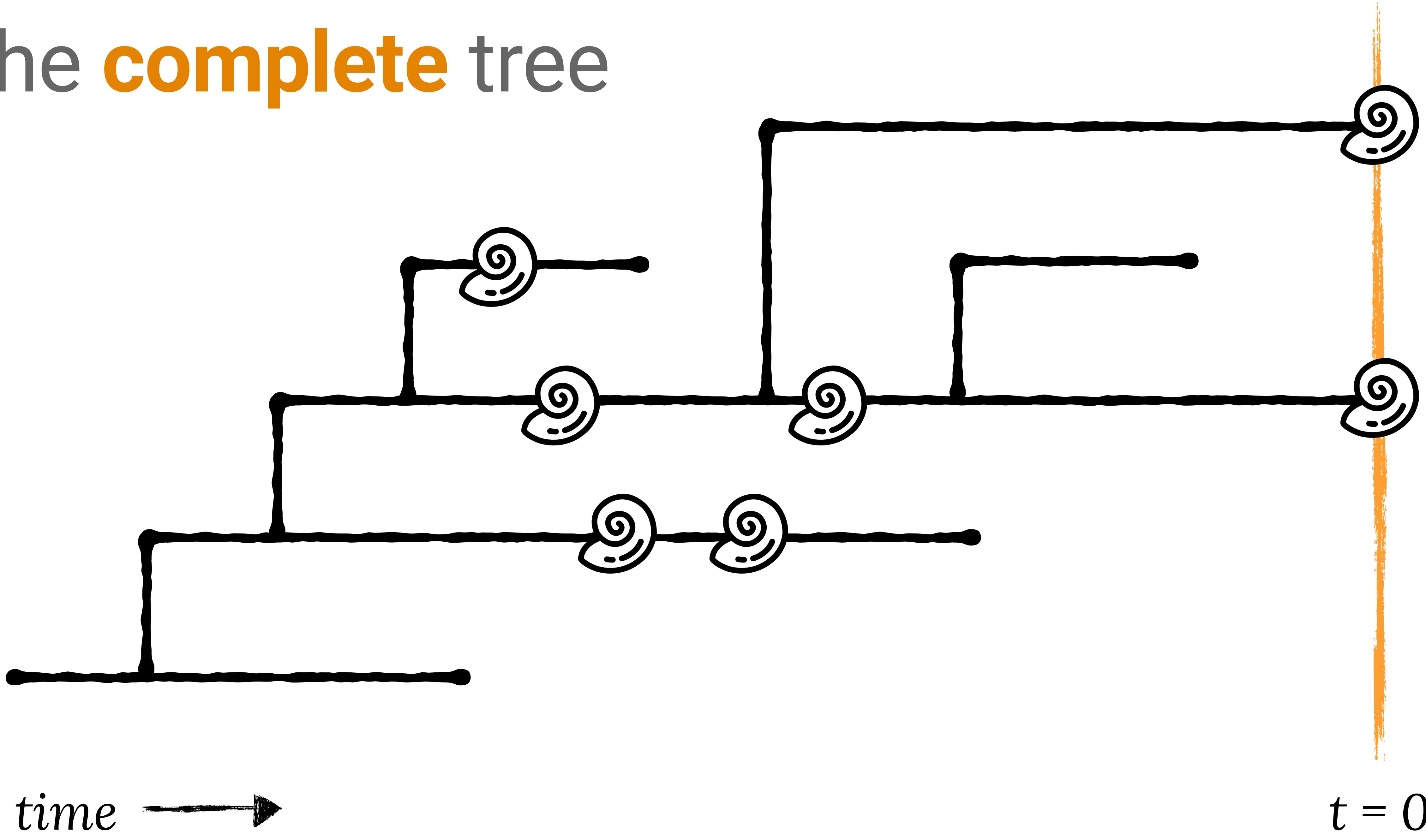


Sampled ancestors

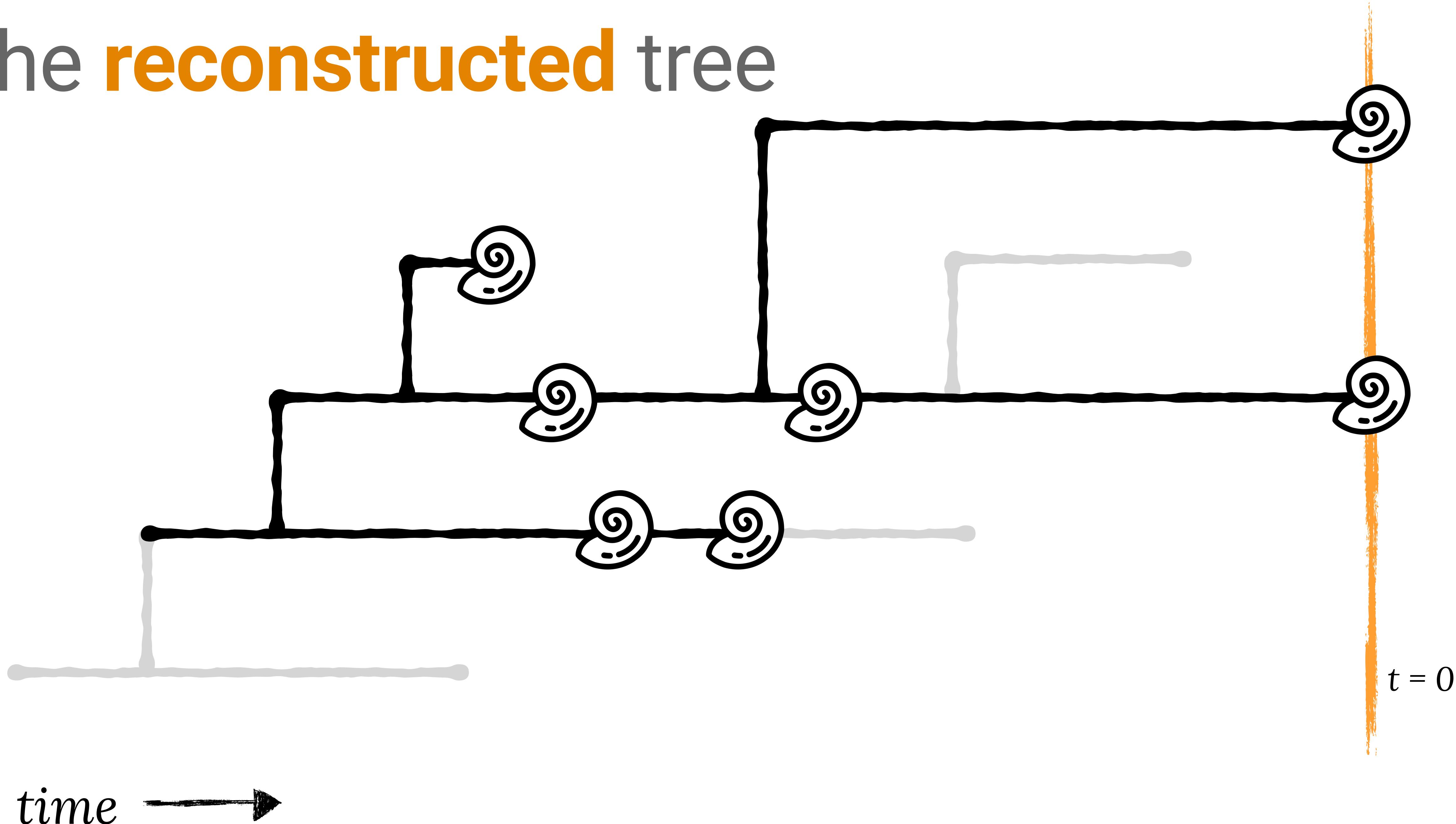
Some lineages go completely unsampled

$t = 0$

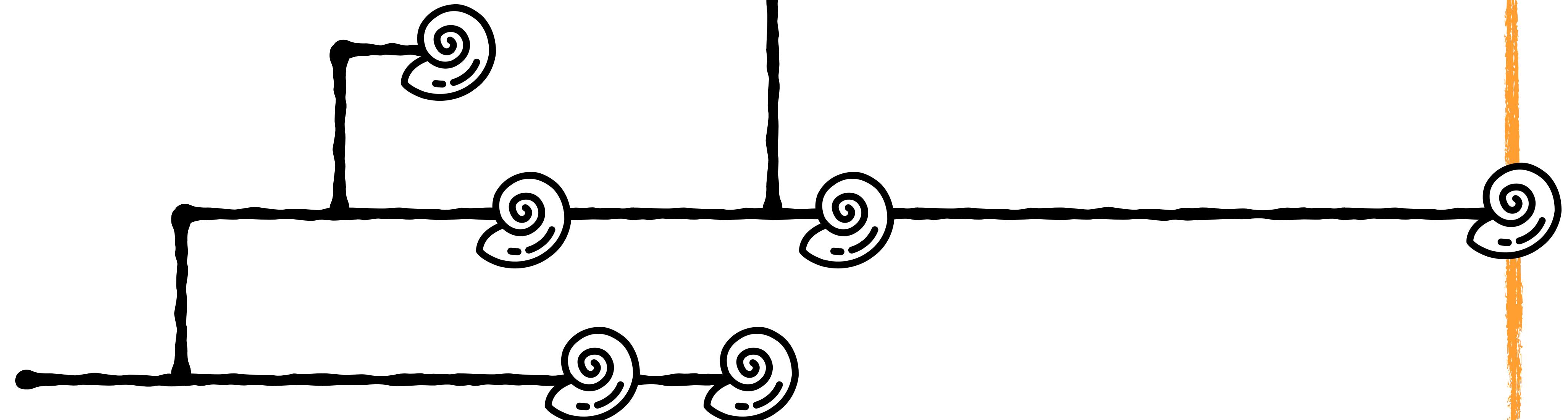
# The complete tree



# The reconstructed tree



The **fossilised birth-death (FBD) process** allows us to calculate the probability of observing the reconstructed tree

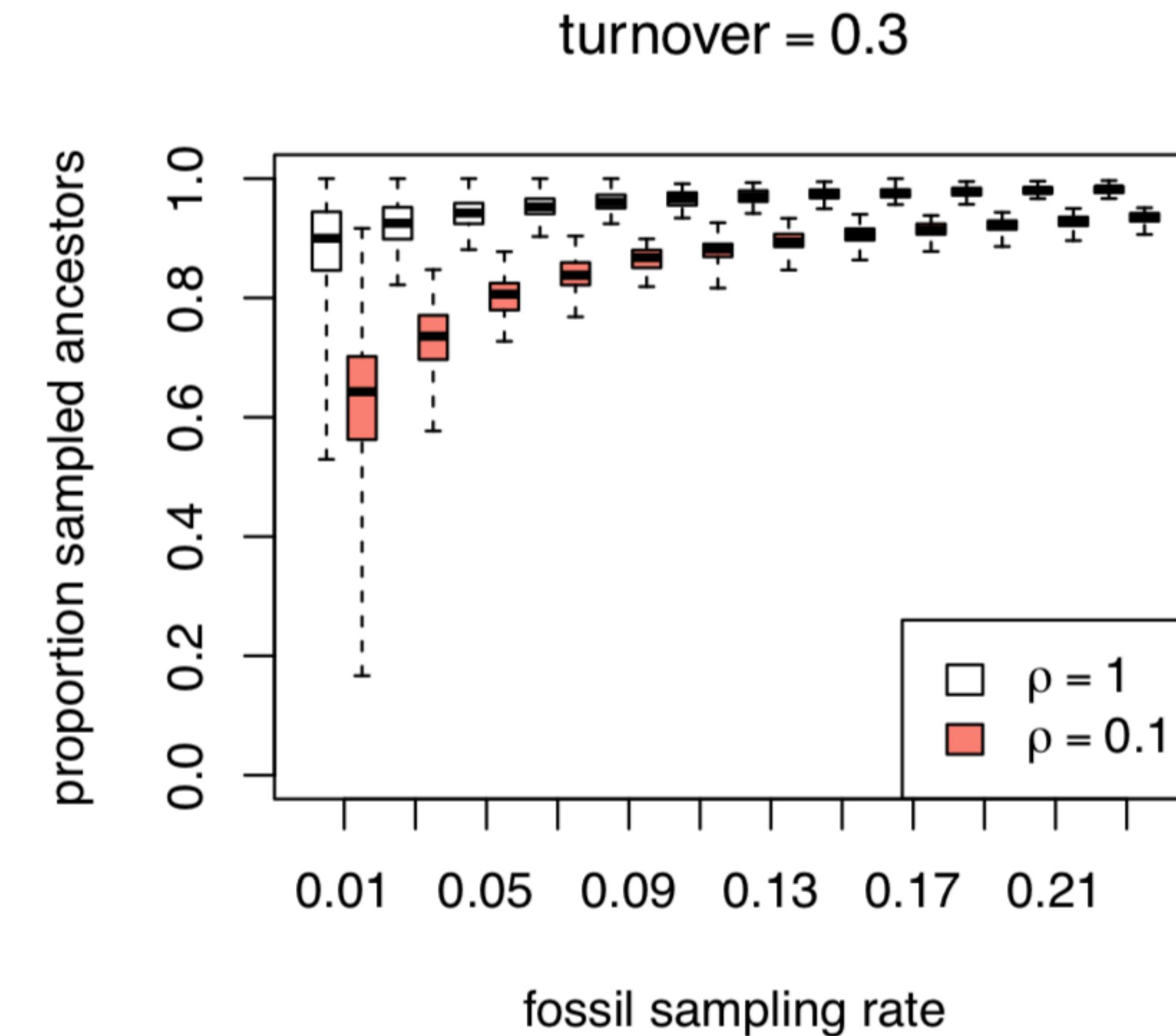
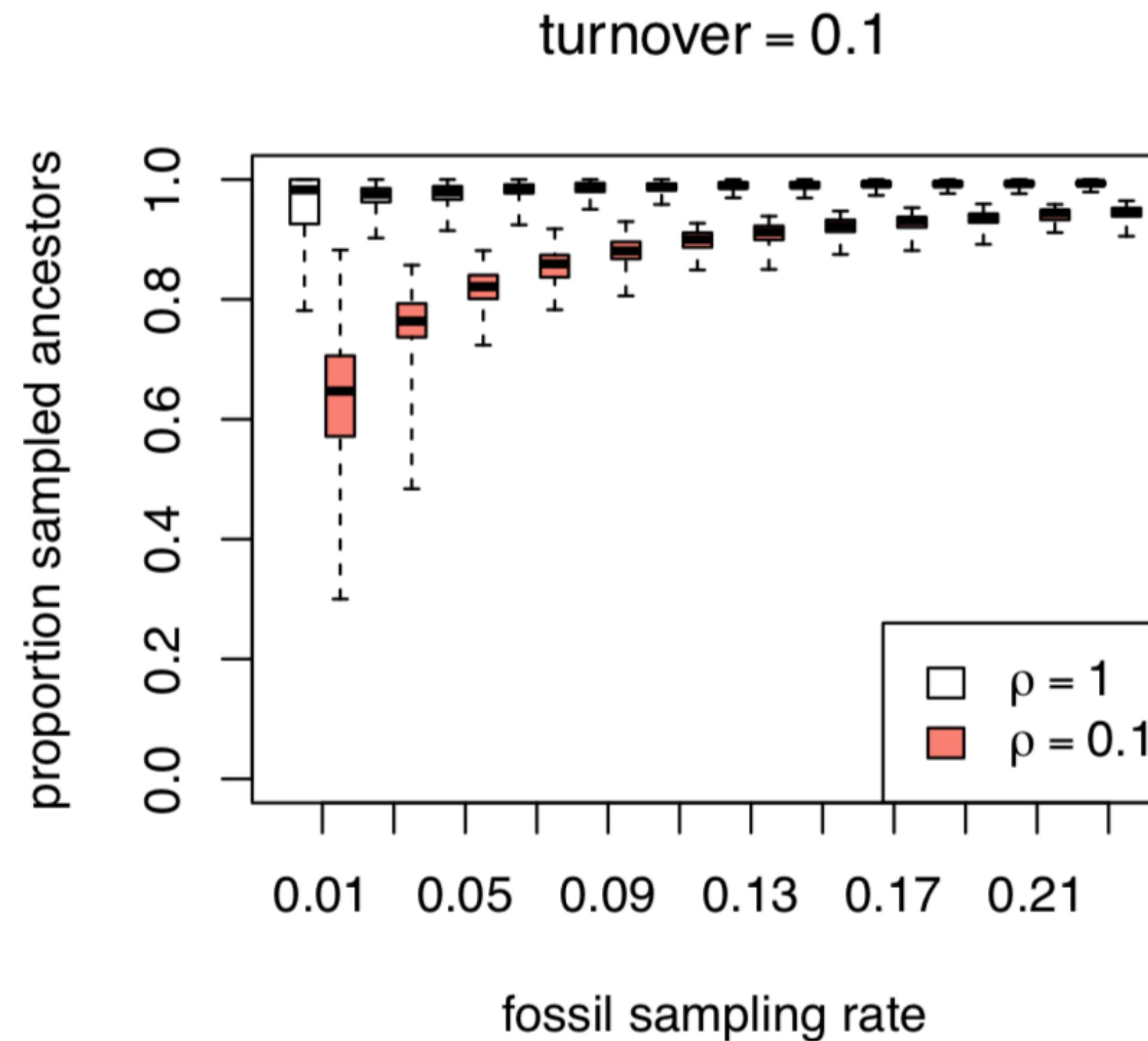


$$P(E | \text{snail}, \lambda, \mu, \psi, \rho)$$

Sampling-through-time in birth-death trees. Stadler. (2010)  
First implemented: Heath et al. (2014) and Gavryushkina et al. (2014)

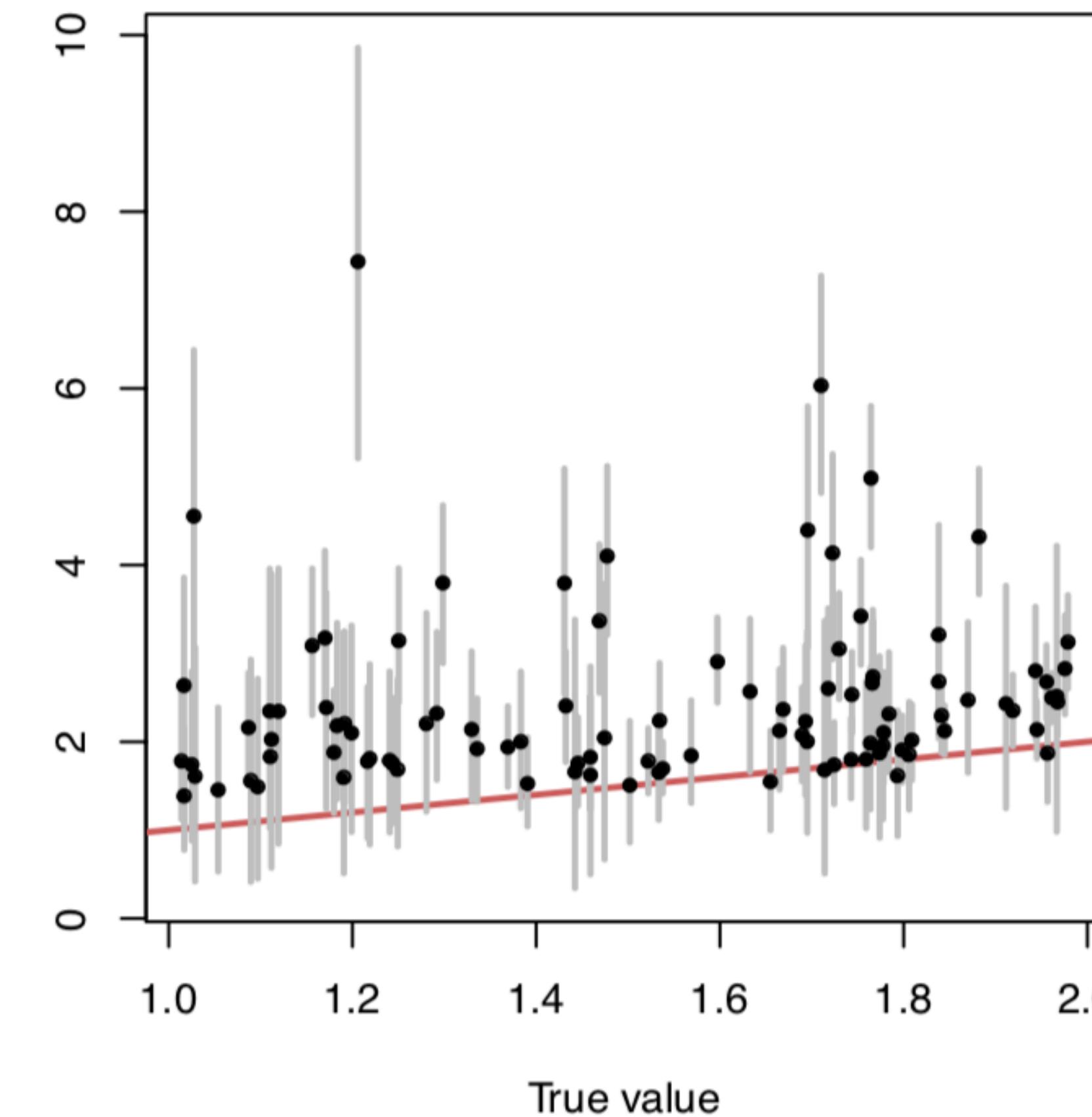
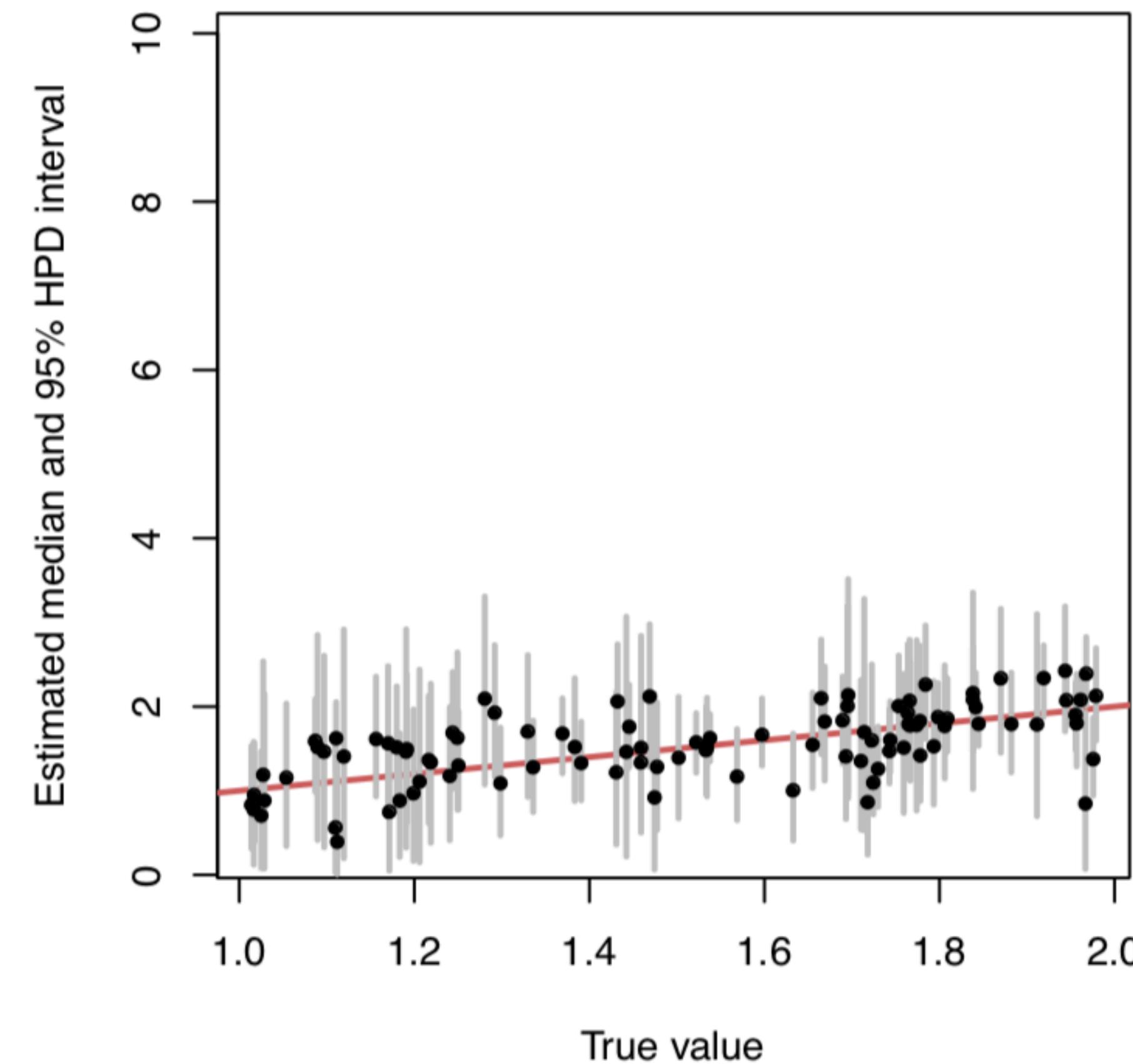
# Sampled ancestors

The proportion increases with higher turnover (birth - death) or higher sampling

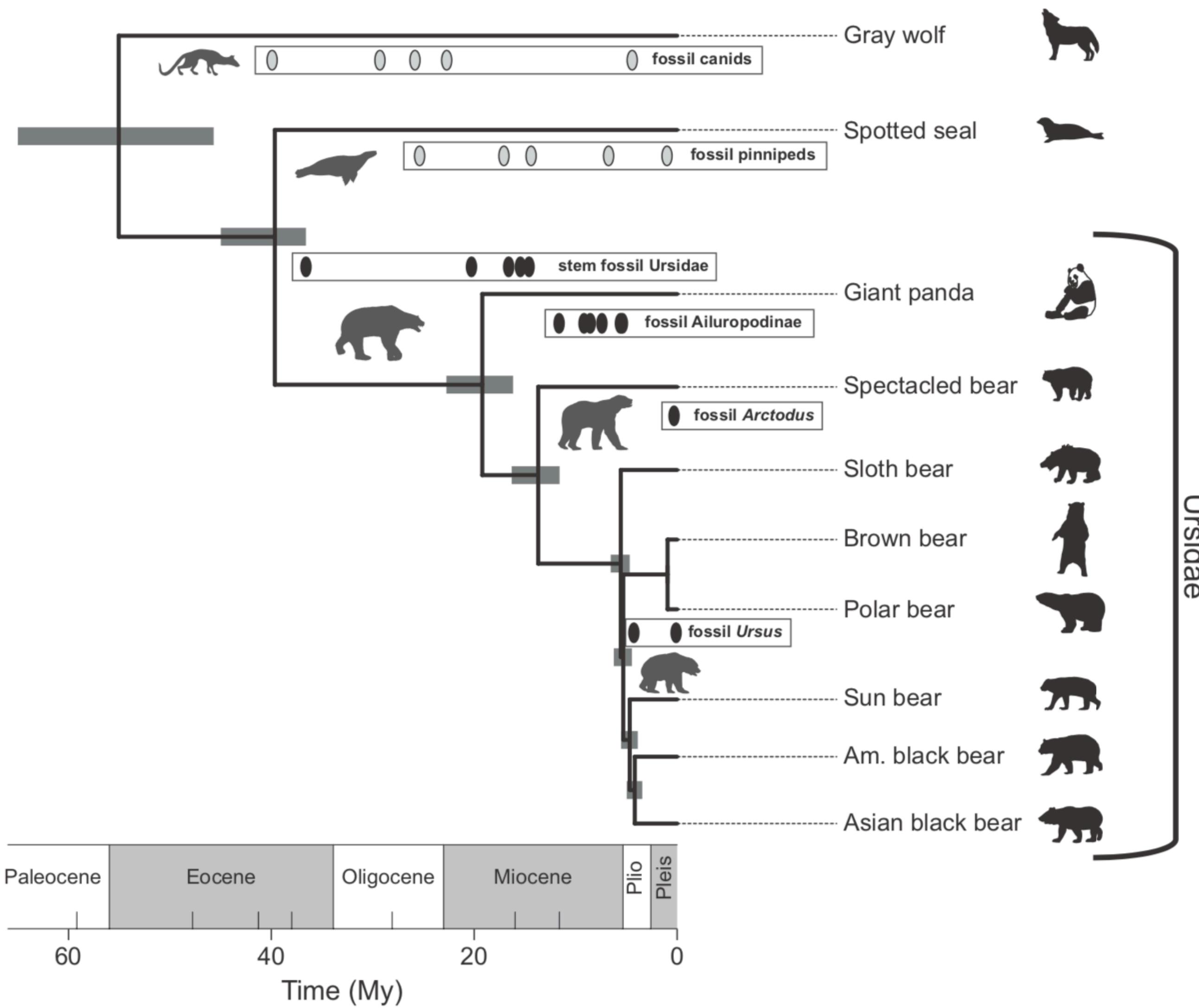


# Sampled ancestors

Ignoring sampled ancestors can lead to inaccurate parameter estimates



# Time calibrated tree of living and fossil bears



First application of the FBD model.

Fossils are incorporated via constraints, not character data. Their precise placement can not be inferred, but this uncertainty will be reflected in the posterior

# Exercise

# Fossils can be incorporated via **taxonomy** or **character data** (total-evidence)

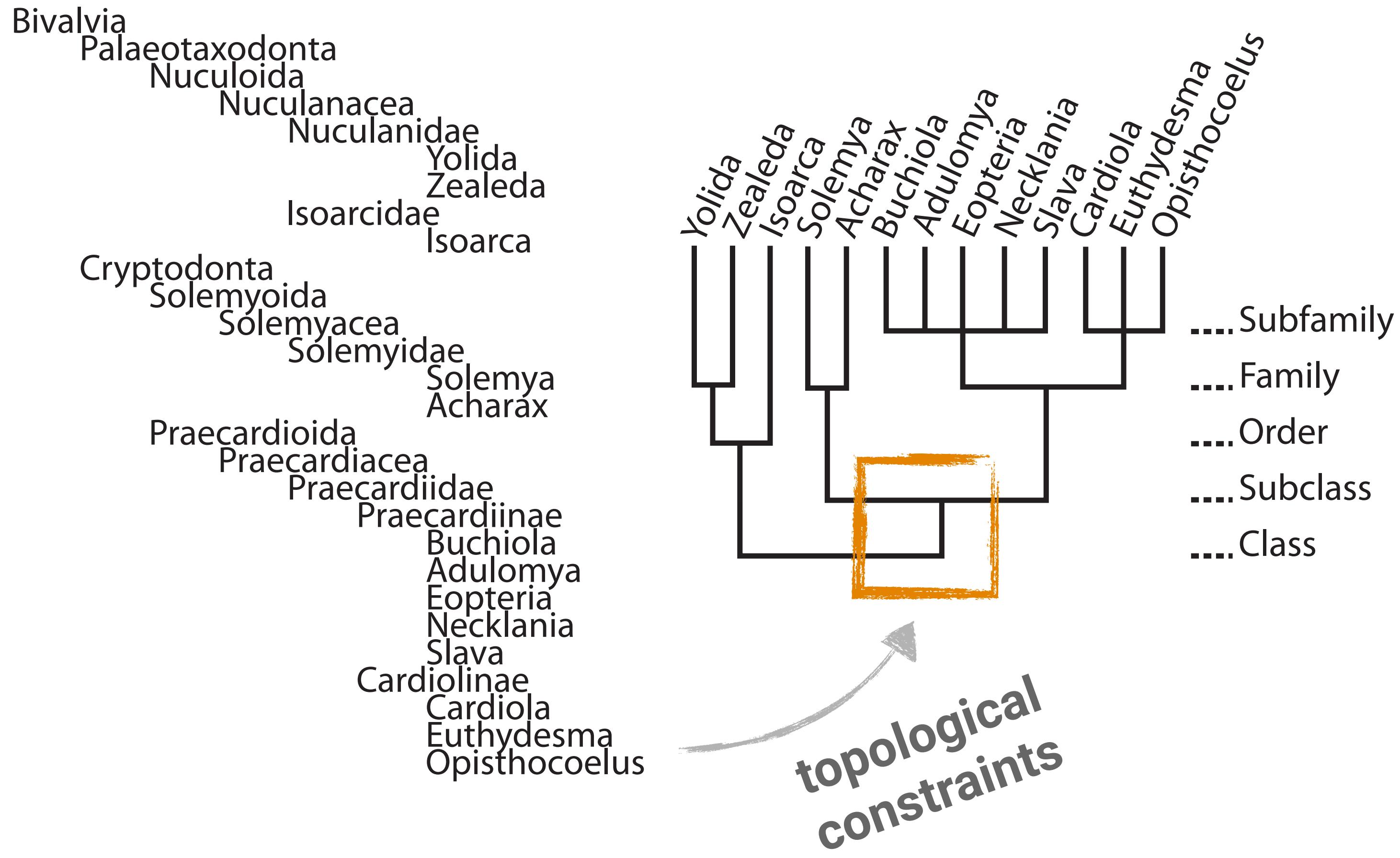
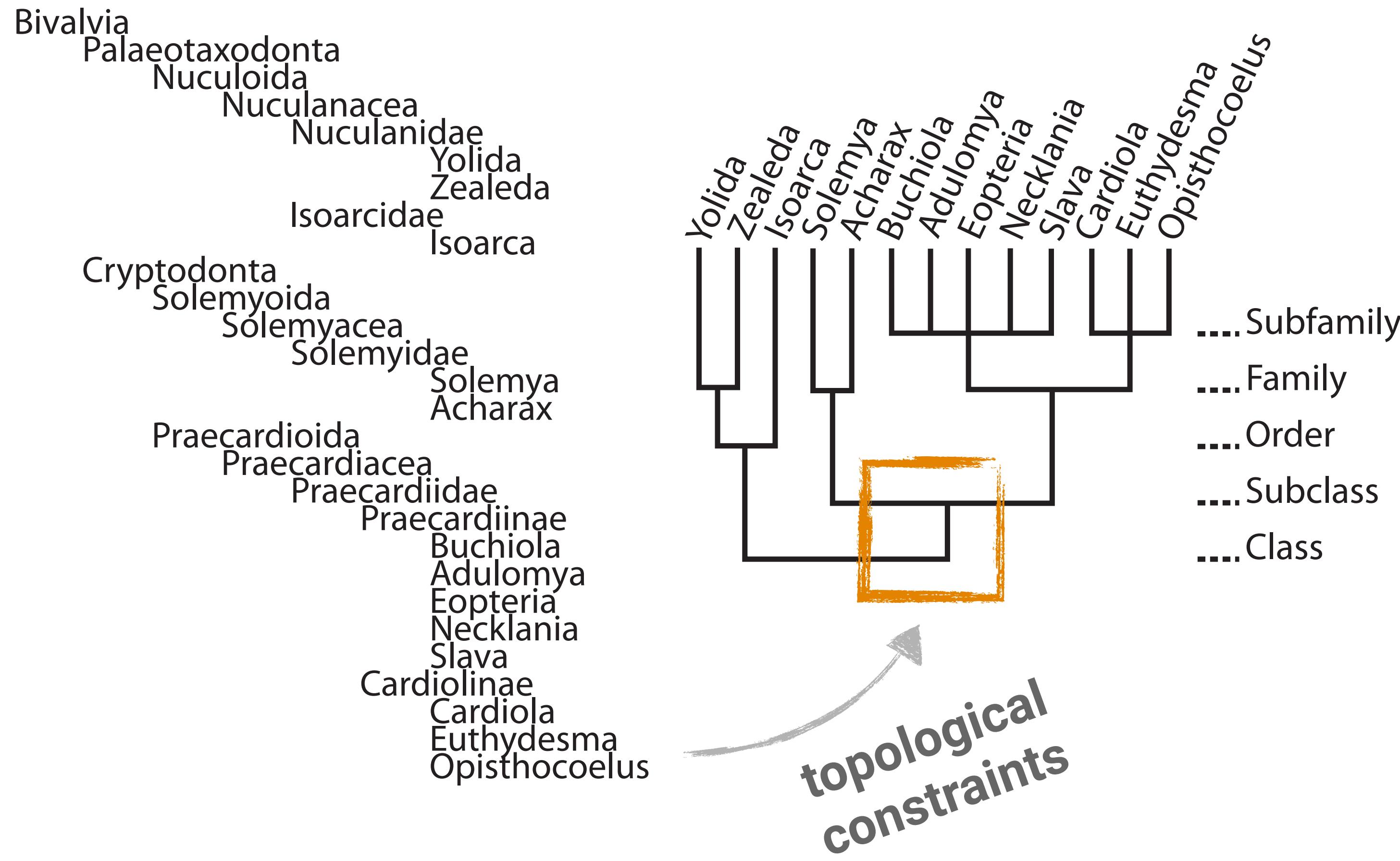


Image source Soul & Friedman (2015)

# Fossils can be incorporated via **taxonomy** or **character data** (total-evidence)



**ATAT...**

**TCACT...**

**?????...**

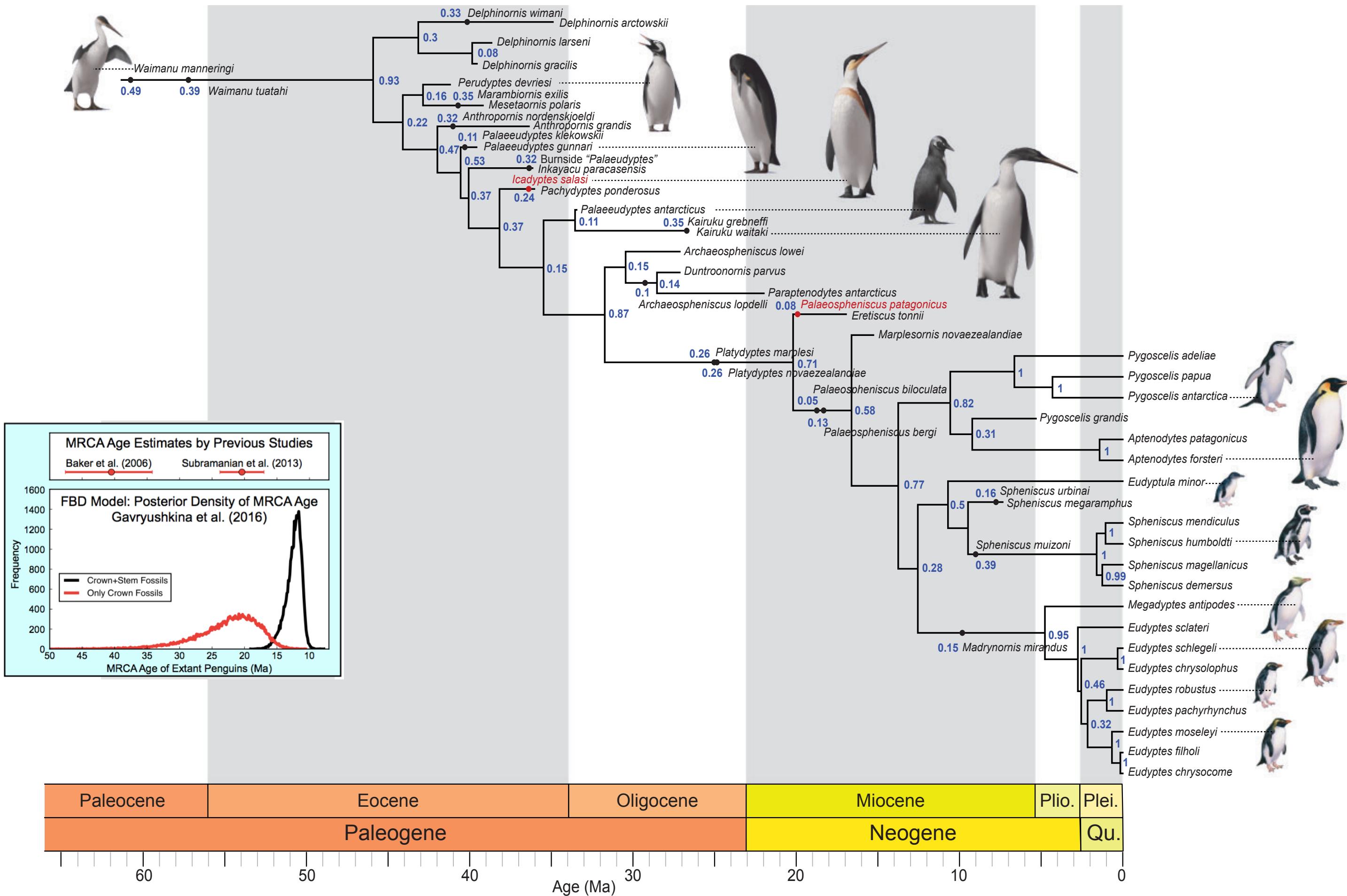
**OR**

**1001...**

**1101...**

**0100...**

# Time calibrated tree of living and fossil penguins

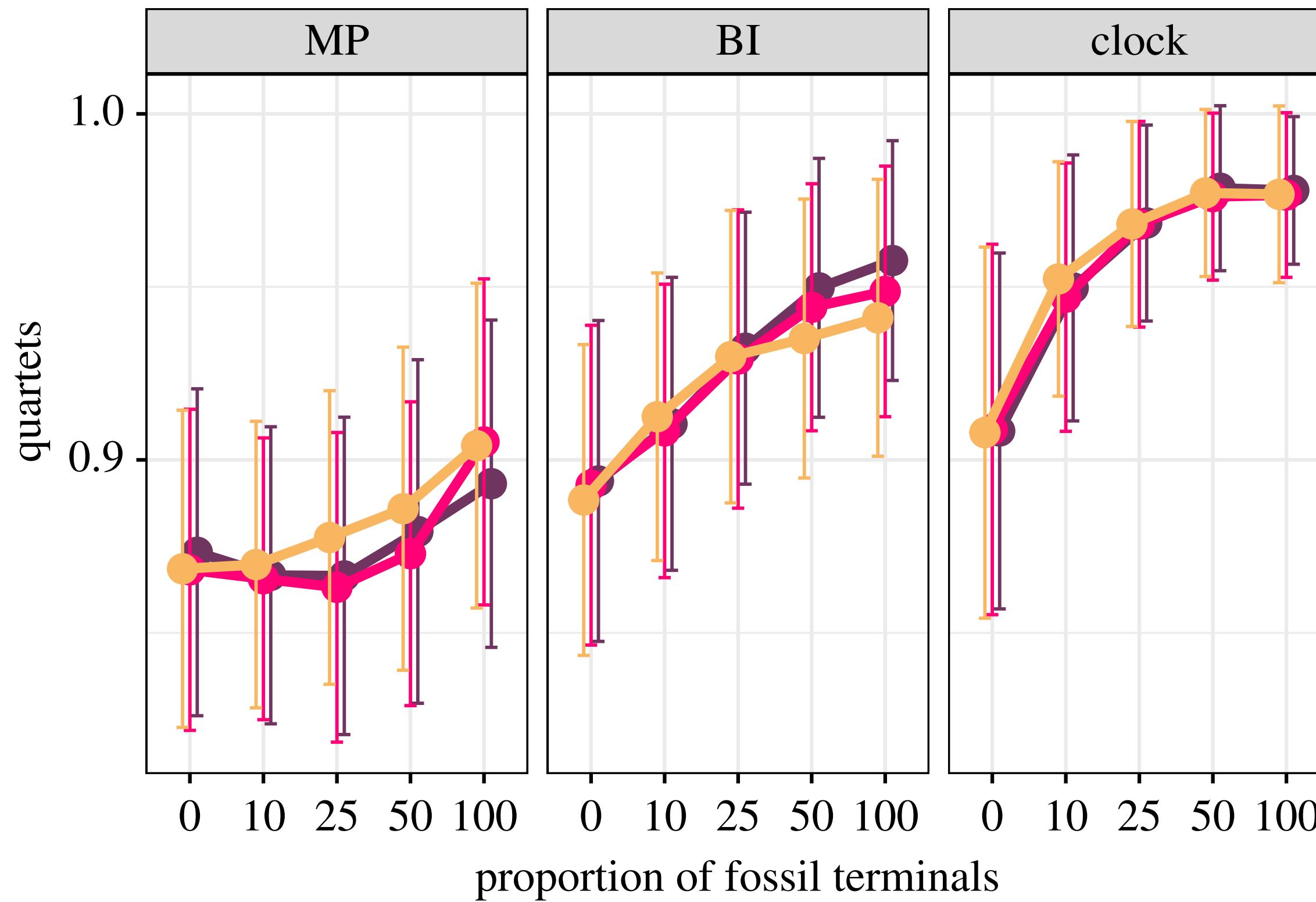


First application of total evidence dating using the FBD model

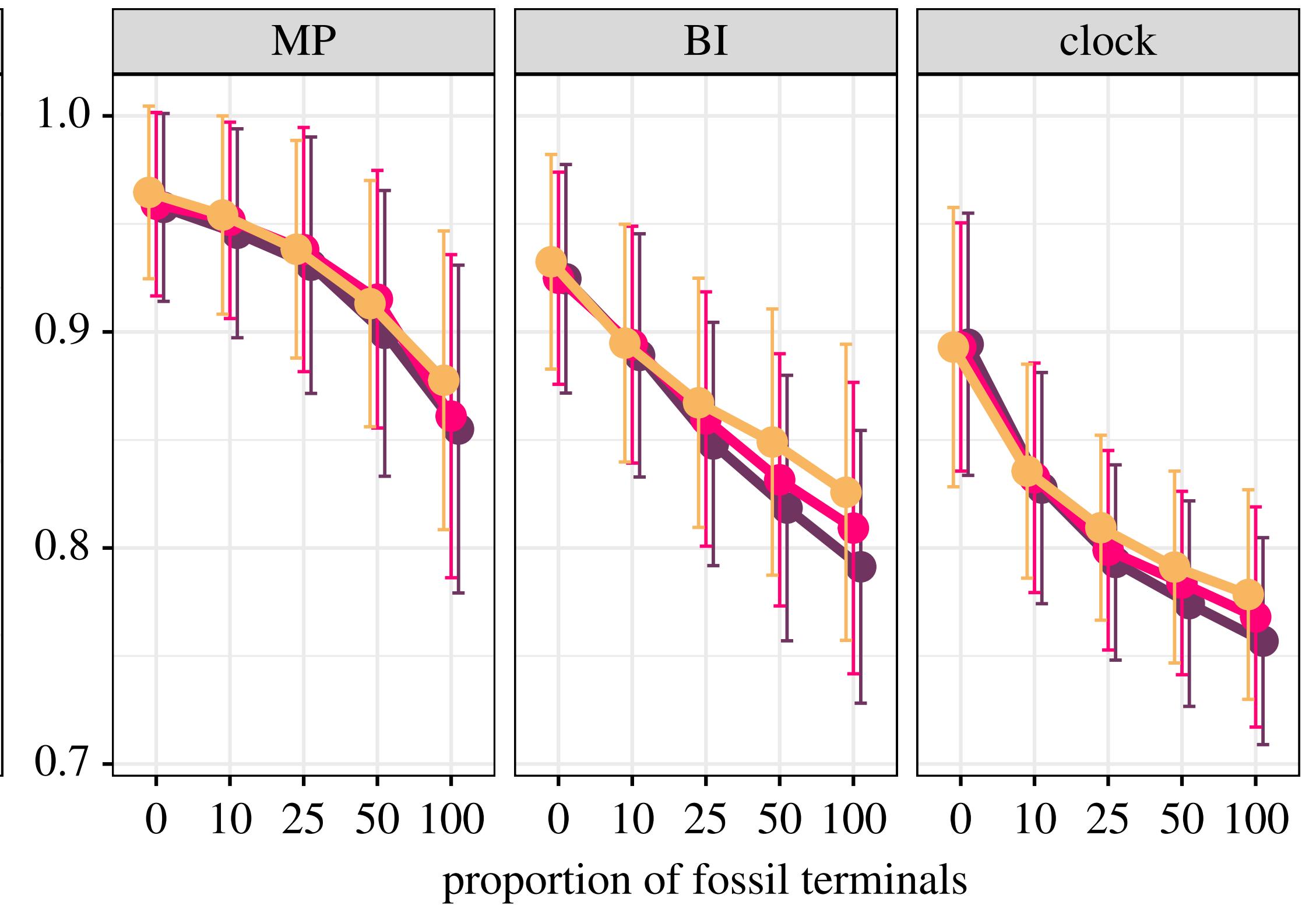
Fossils are incorporated using character data

Gavryushkina et al. (2016)

# accuracy



# precision



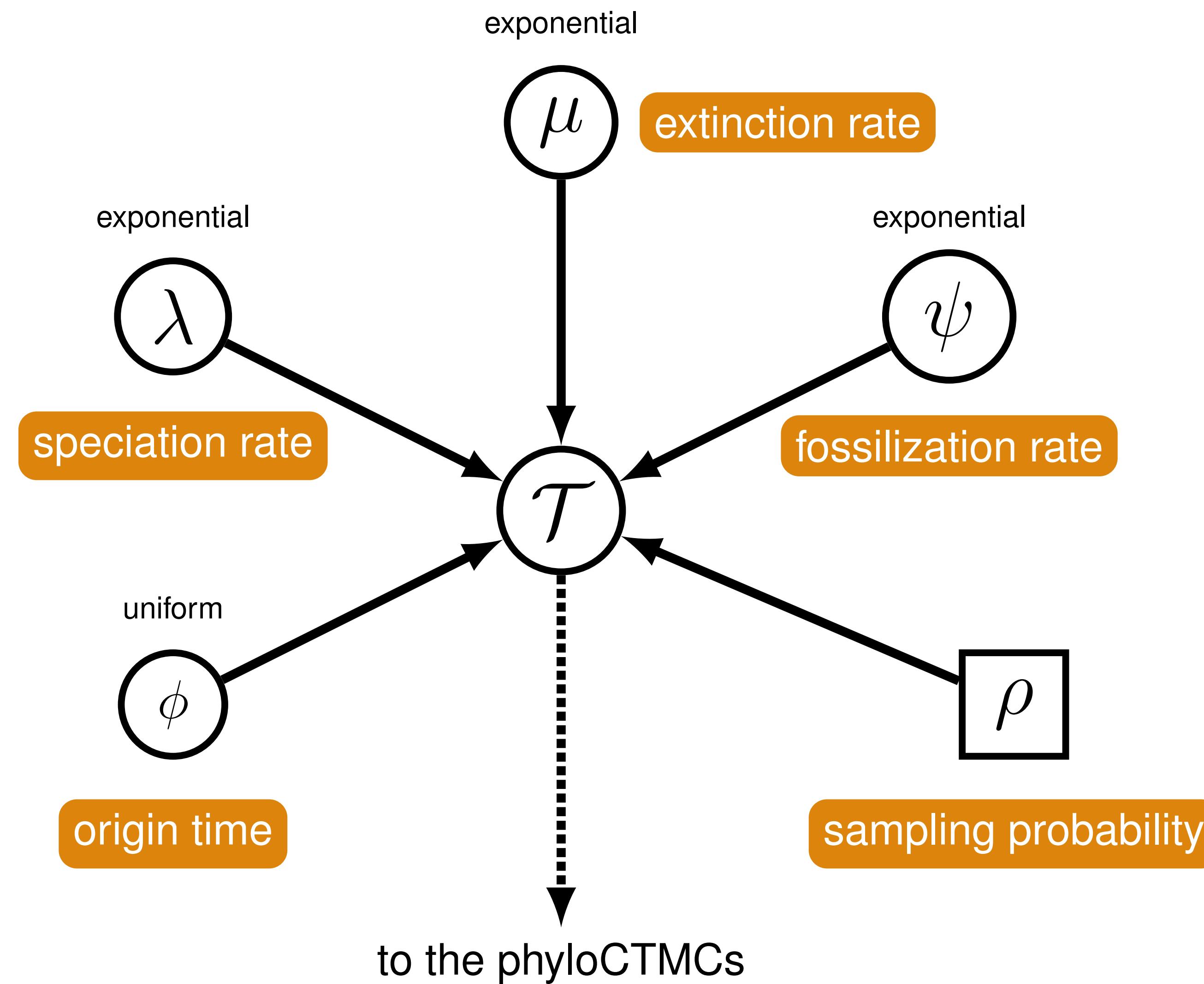
level of missing data      ● none    ● low    ● high

*Fossils improve phylogenetic analyses of morphological characters*  
Koch, Garwood, Parry. 2020. Proc B

# Some notes

- The topology of extant taxa is largely unaffected by how fossils are incorporated
- Fossils *and* age information help inform topology
- Divergence times are much more sensitive to errors in fossil placement and model misspecification
- Total-evidence dating is more robust to model misspecification

# Graphical model representation of the FBDP



# Relationship to (some) other birth-death process models

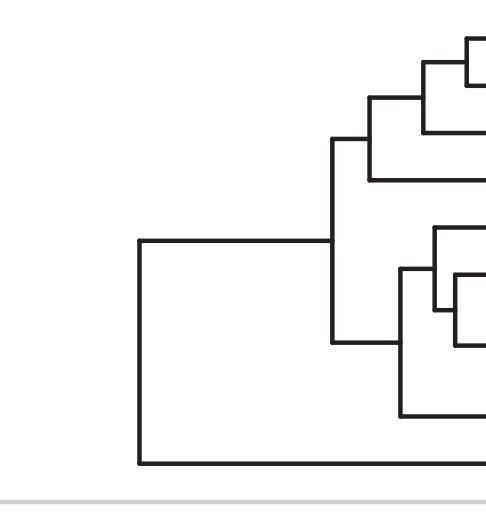
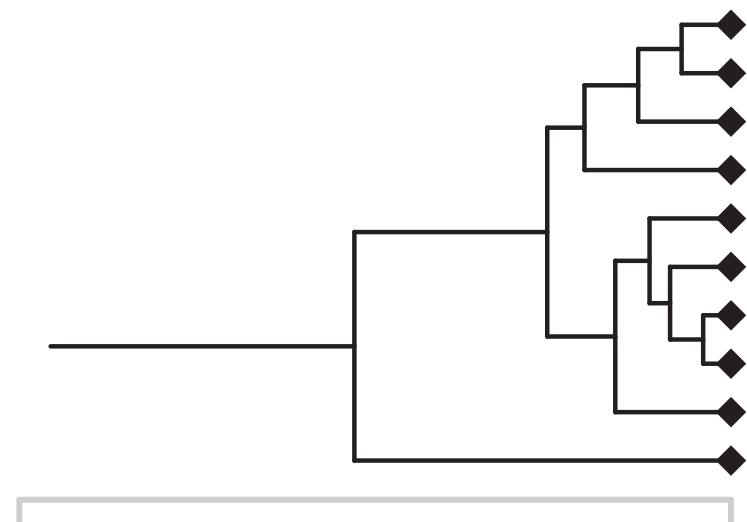
These models are special cases of the FBD process, with fossil sampling ( $\psi$ ) = zero.

We can also use  $\rho$  at  $t > 0$  to model serial sampling.

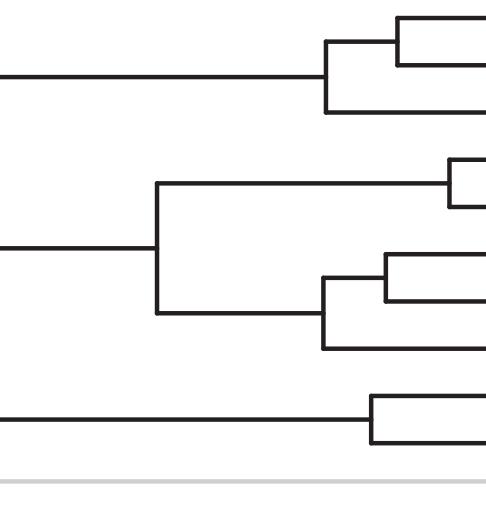
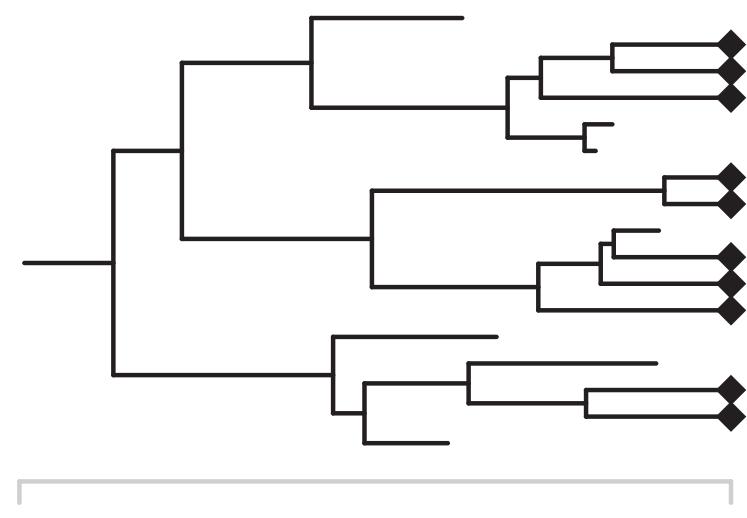
Stadler et al. 2012

See also: Stadler and Yang 2013

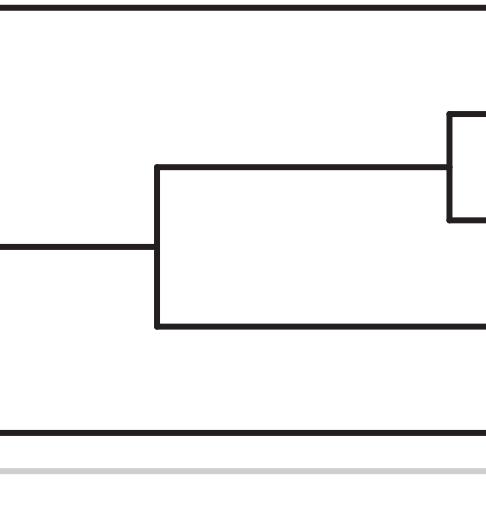
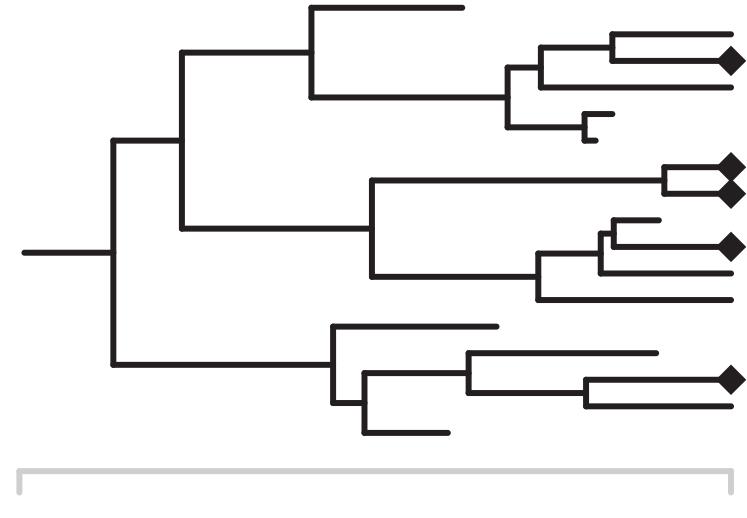
complete vs. reconstructed trees



$$\lambda = 0.1$$

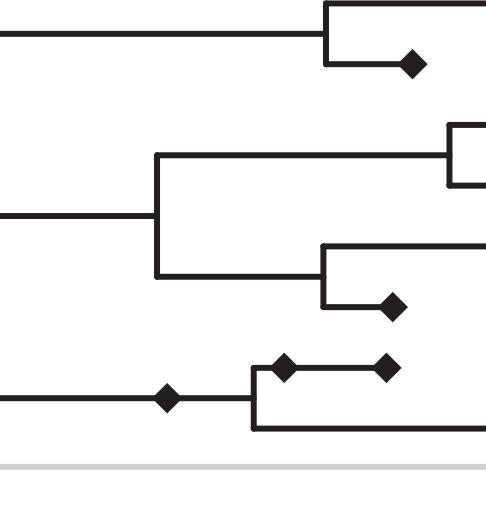


$$\lambda = 0.1 \\ \mu = 0.05$$



$$\lambda = 0.1 \\ \mu = 0.05 \\ \rho = 0.6$$

Yang and  
Rannala  
1997  
Stadler  
2009



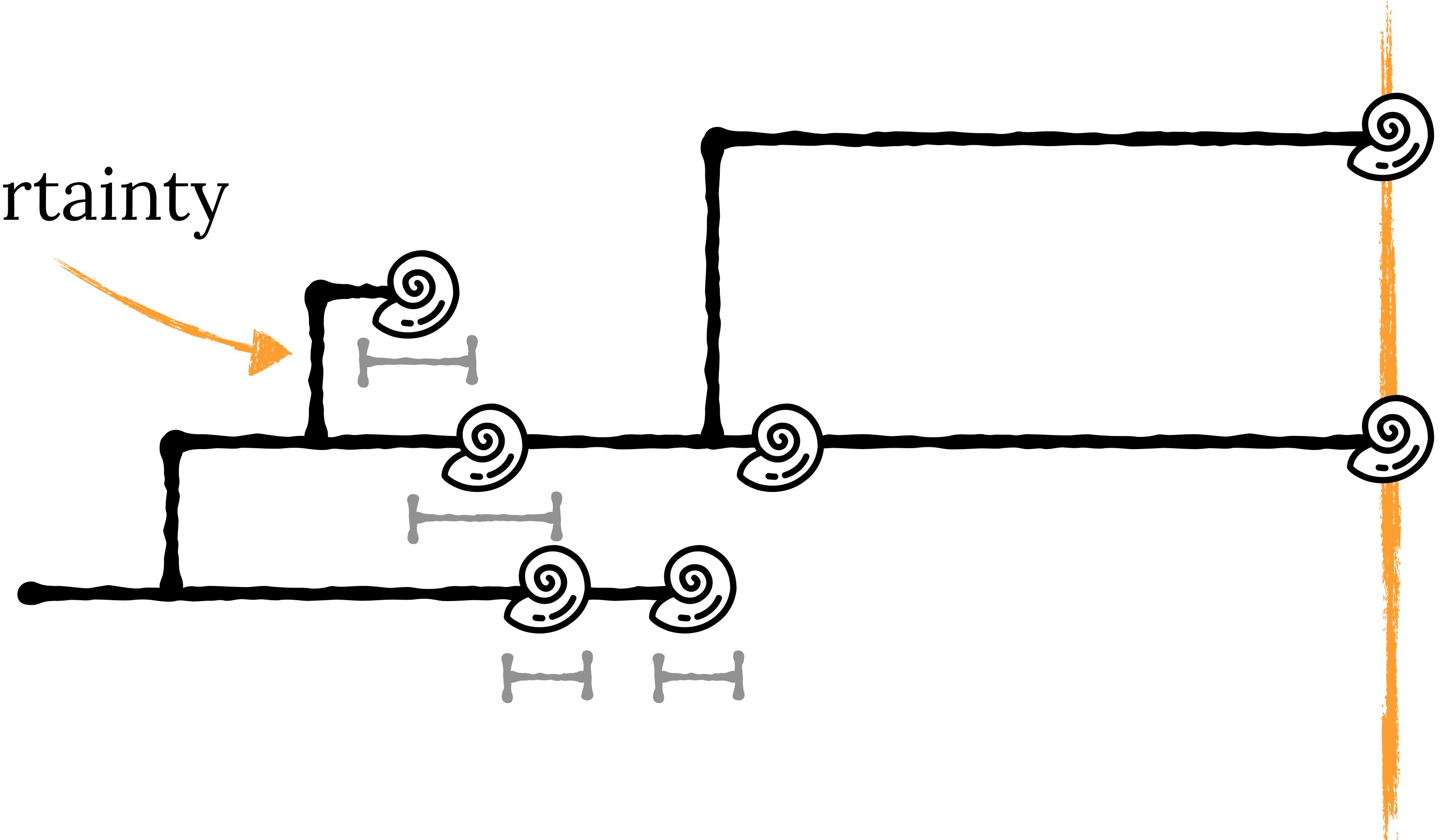
$$\lambda = 0.1 \\ \mu = 0.05 \\ \rho = 0.6 \\ \psi = 0.05$$

Stadler  
2010

# Sample age uncertainty



age uncertainty

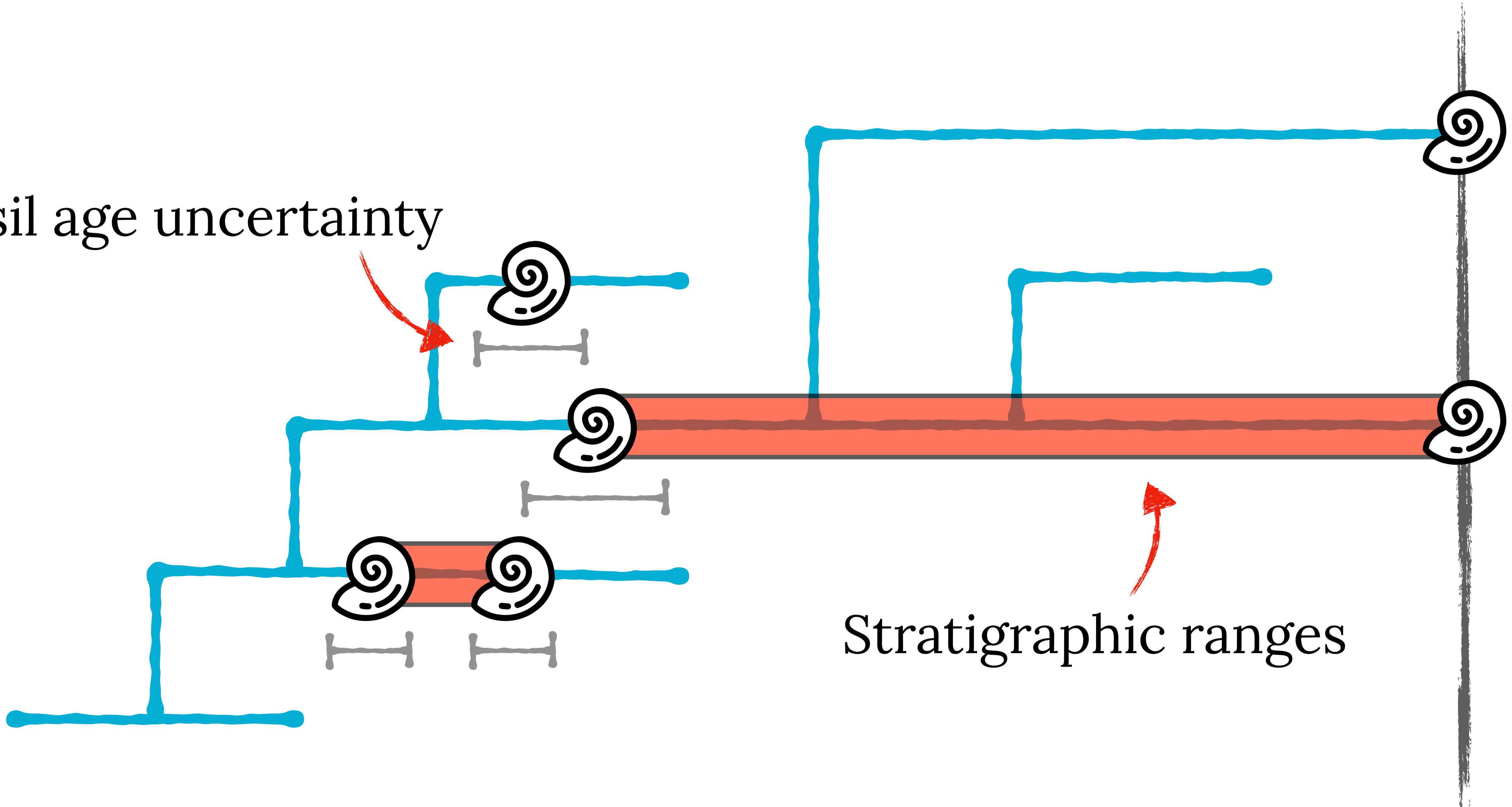


*Ignoring Fossil Age Uncertainty Leads to Inaccurate Topology in Time Calibrated Tree Inference*

Barido-Sottani et al. 2018, 2020

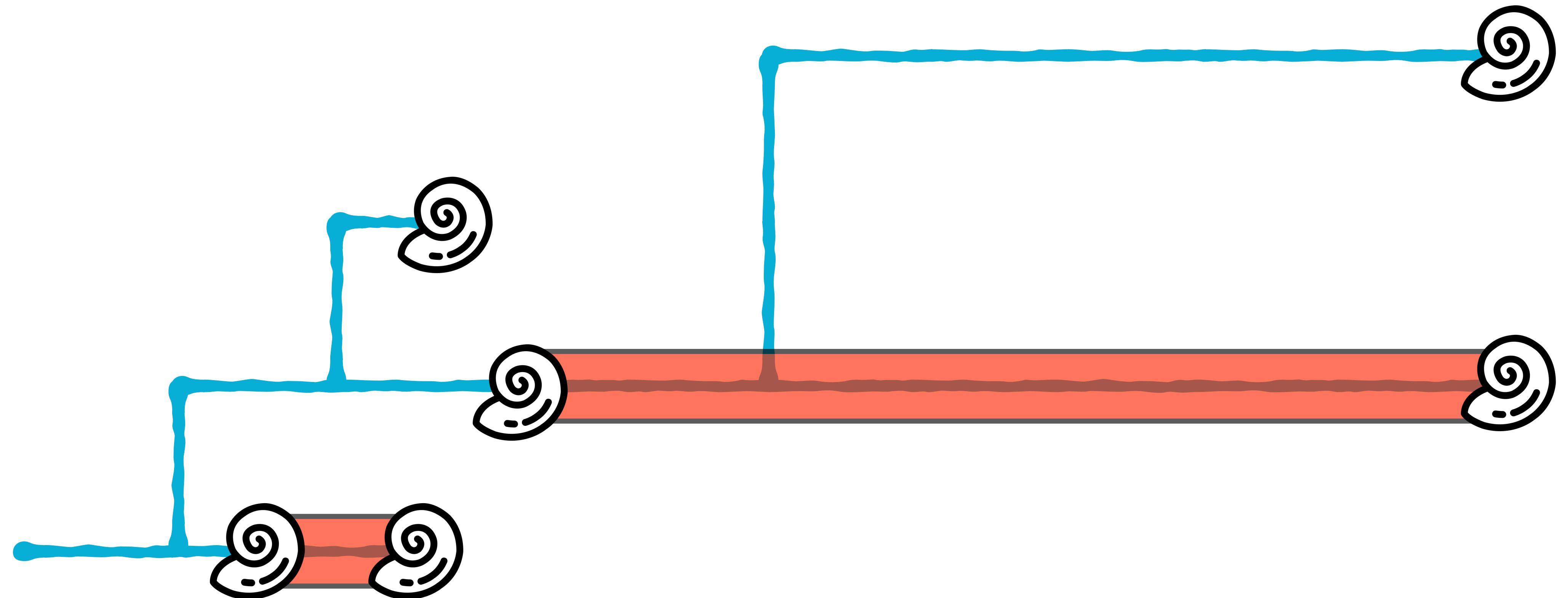
*Putting the F in FBD analyses: tree constraints or morphological data?* Barido-Sottani et al. 2023<sub>43</sub>

fossil age uncertainty

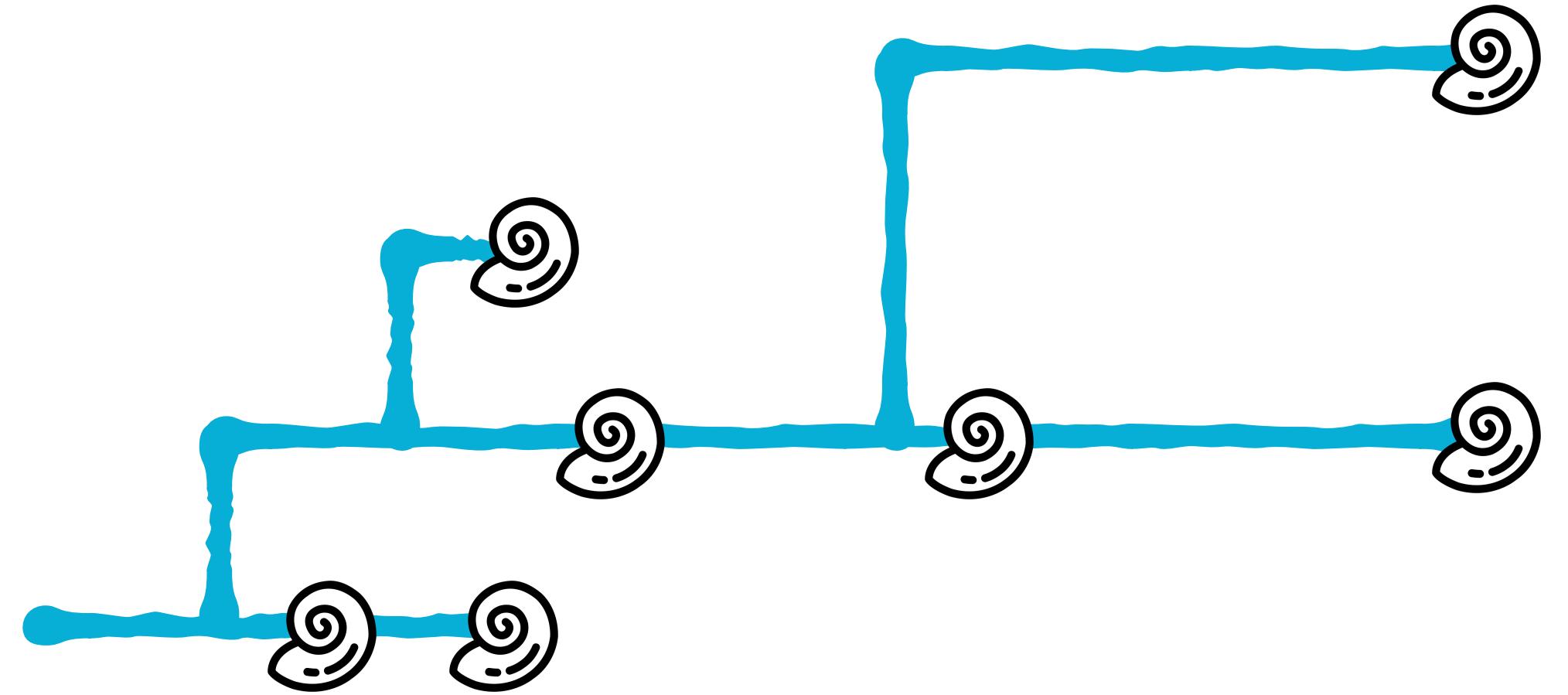


*The inseparability of sampling and time and its influence on attempts to unify the molecular & fossil records*  
Hopkins et al. 2018. Paleobiology

# The fossilised birth-death range process



*The fossilised birth-death model for the analysis of stratigraphic range data under different speciation modes  
Stadler et al. 2018. JTB*

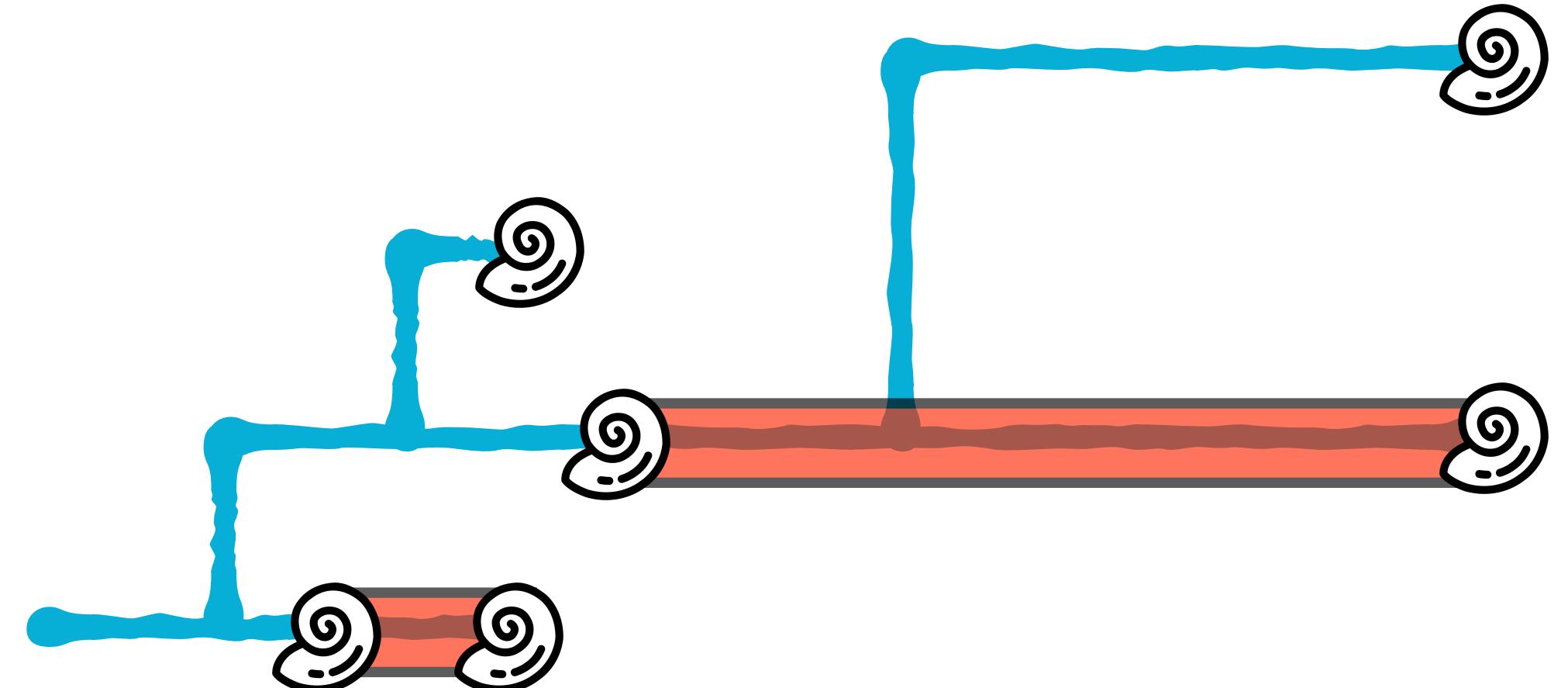


FBD process for analysis of specimen level data

$$P(E | \text{spiral}^{\lambda, \mu, \psi}, p)$$

The FBD range process for analysis of stratigraphic ranges

$$\Pr(E | \text{spiral}^{\lambda, \mu, \psi})$$



# Exercise

# Phylogenetics

Diversification rate estimation

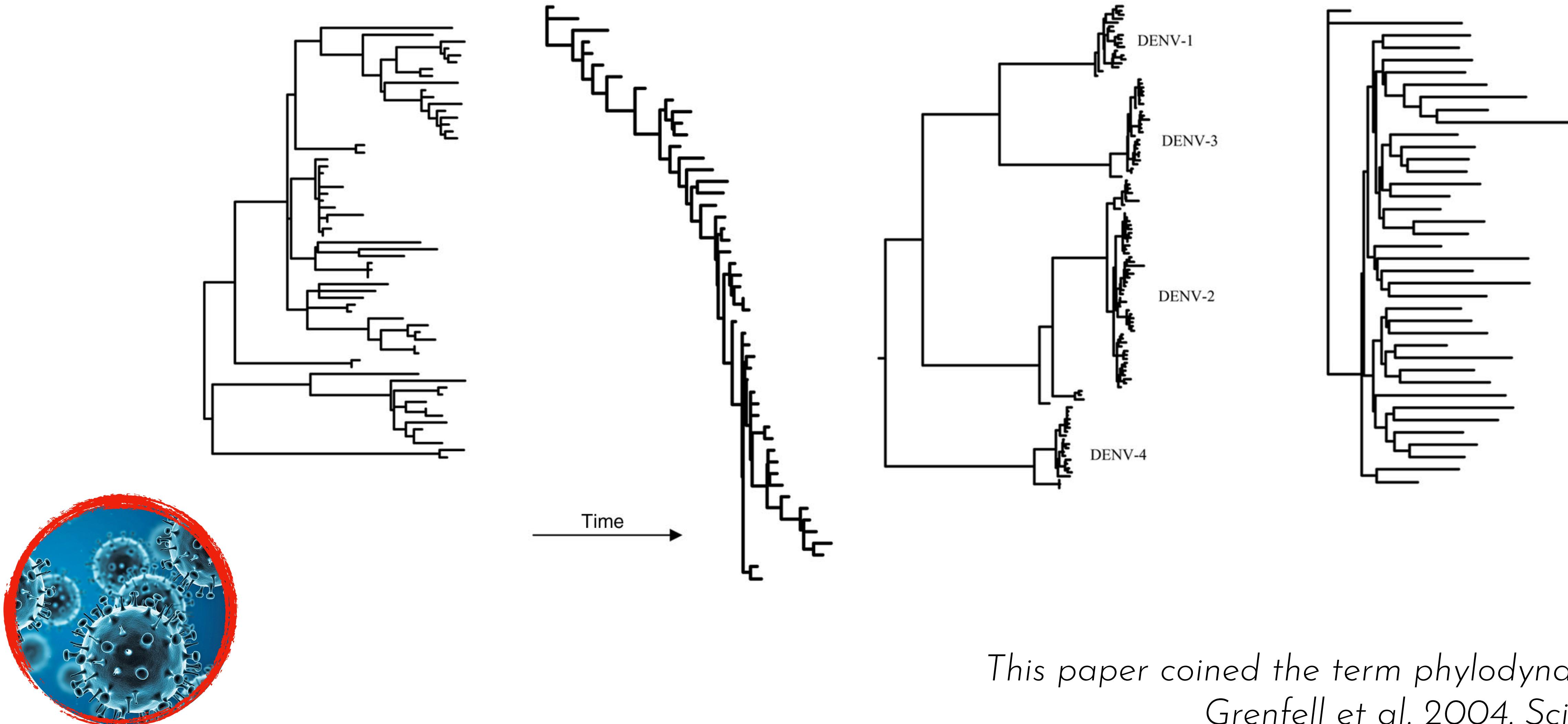
# Bayesian divergence time estimation

$$P(E \mid \lambda, \mu, \psi, p, O, t \mid 0101\dots, 1101\dots, 0100\dots, \text{snail}) =$$

probability of the  
time tree

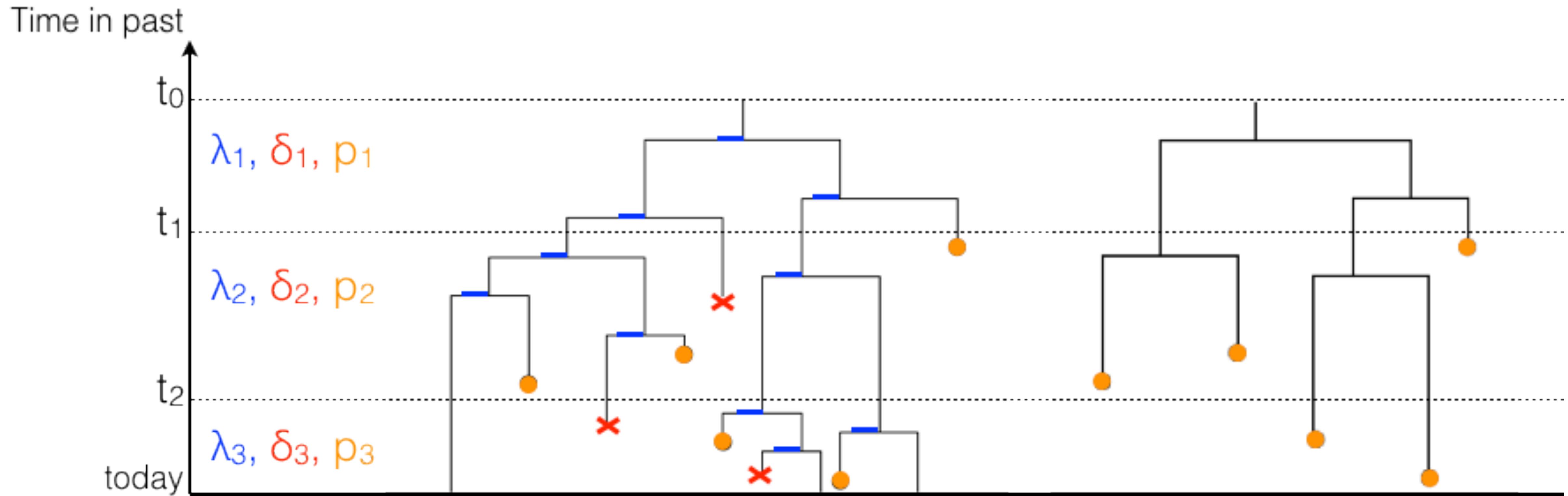
$$\frac{P(0101\dots, 1101\dots, 0100\dots \mid E) P(E \mid \lambda, \mu, \psi, p, O, t) P(\lambda, \mu, \psi, p) P(O) P(t)}{P(0101\dots, 1101\dots, 0100\dots, \text{snail})}$$

# Tree shape is informative about underlying dynamics

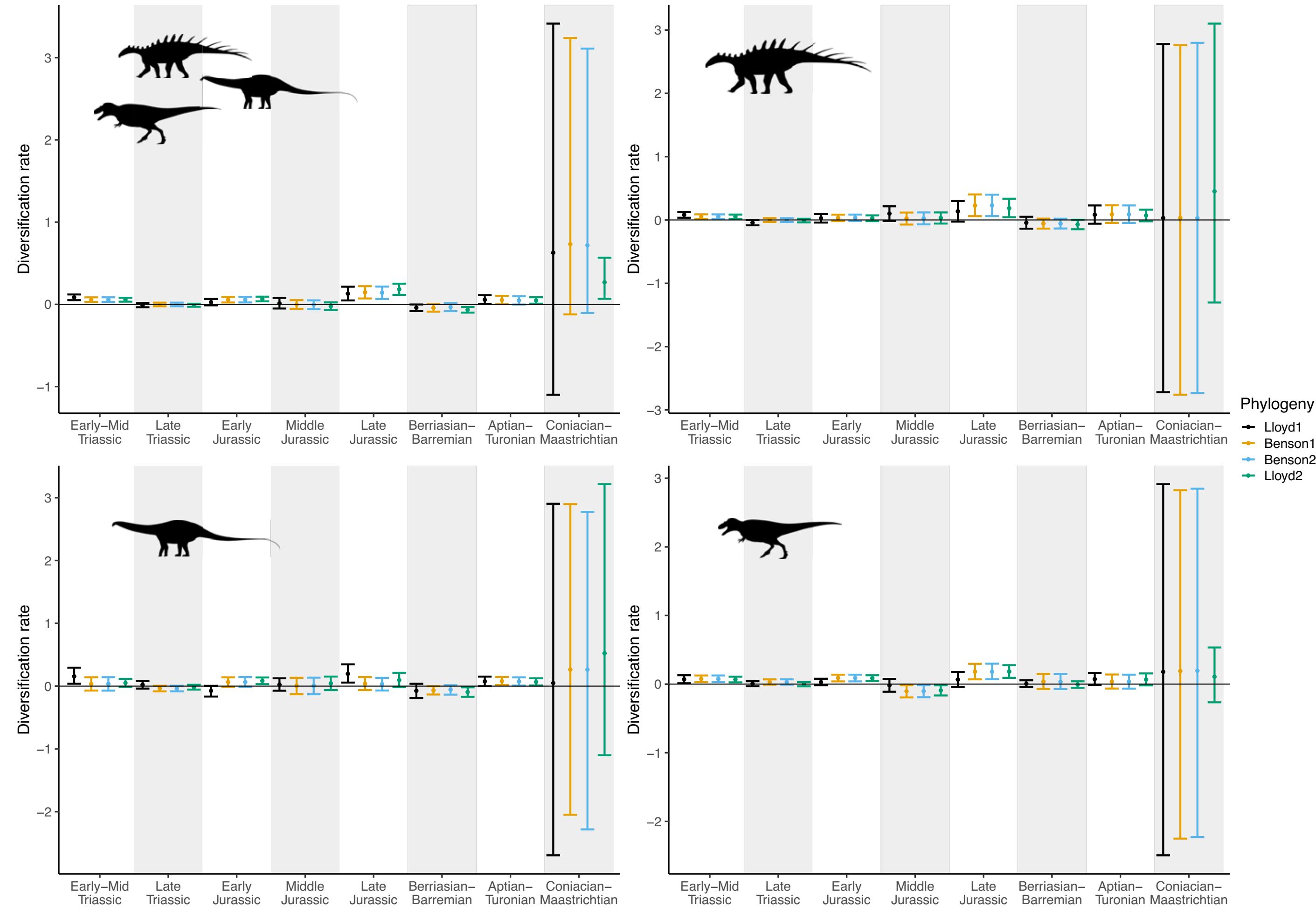


# The skyline birth-death process

First used for tracking the spread of infectious diseases



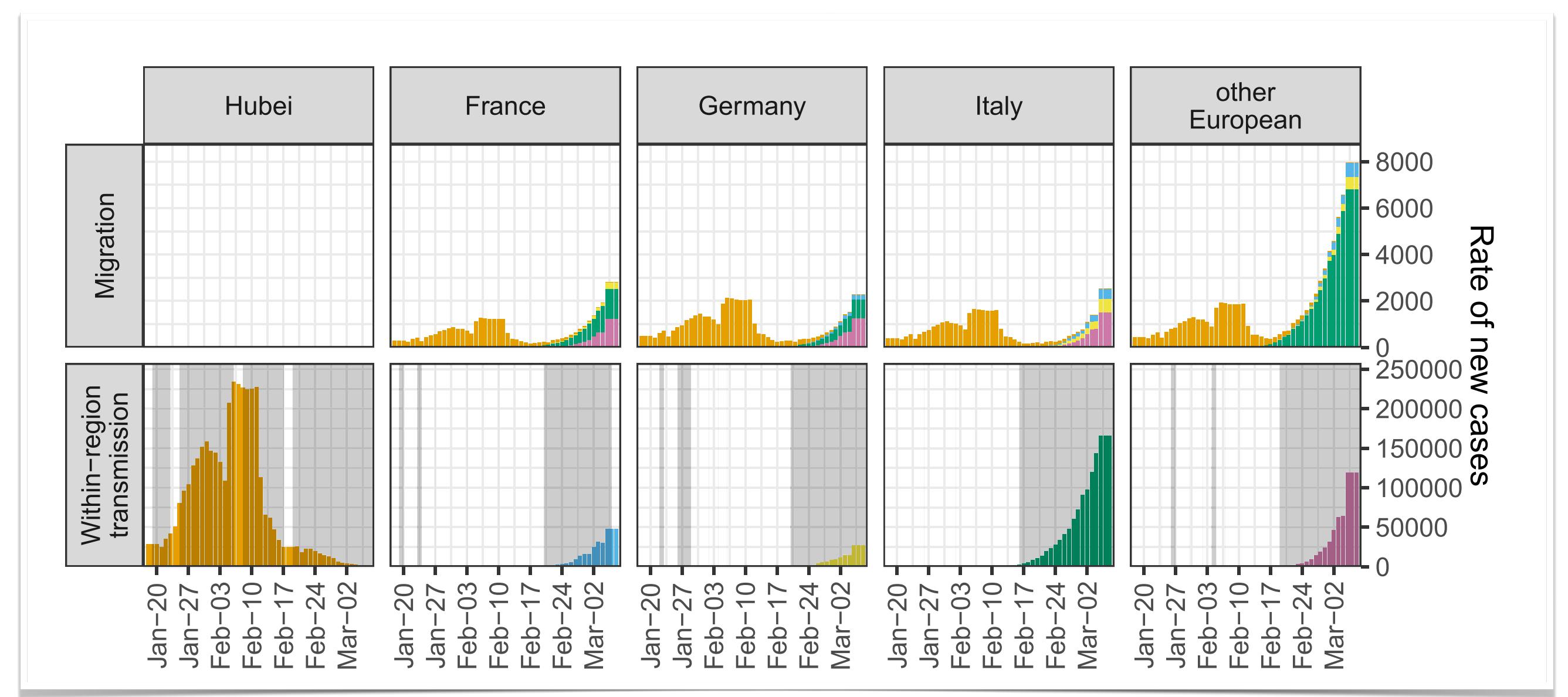
# Macroevolutionary case study



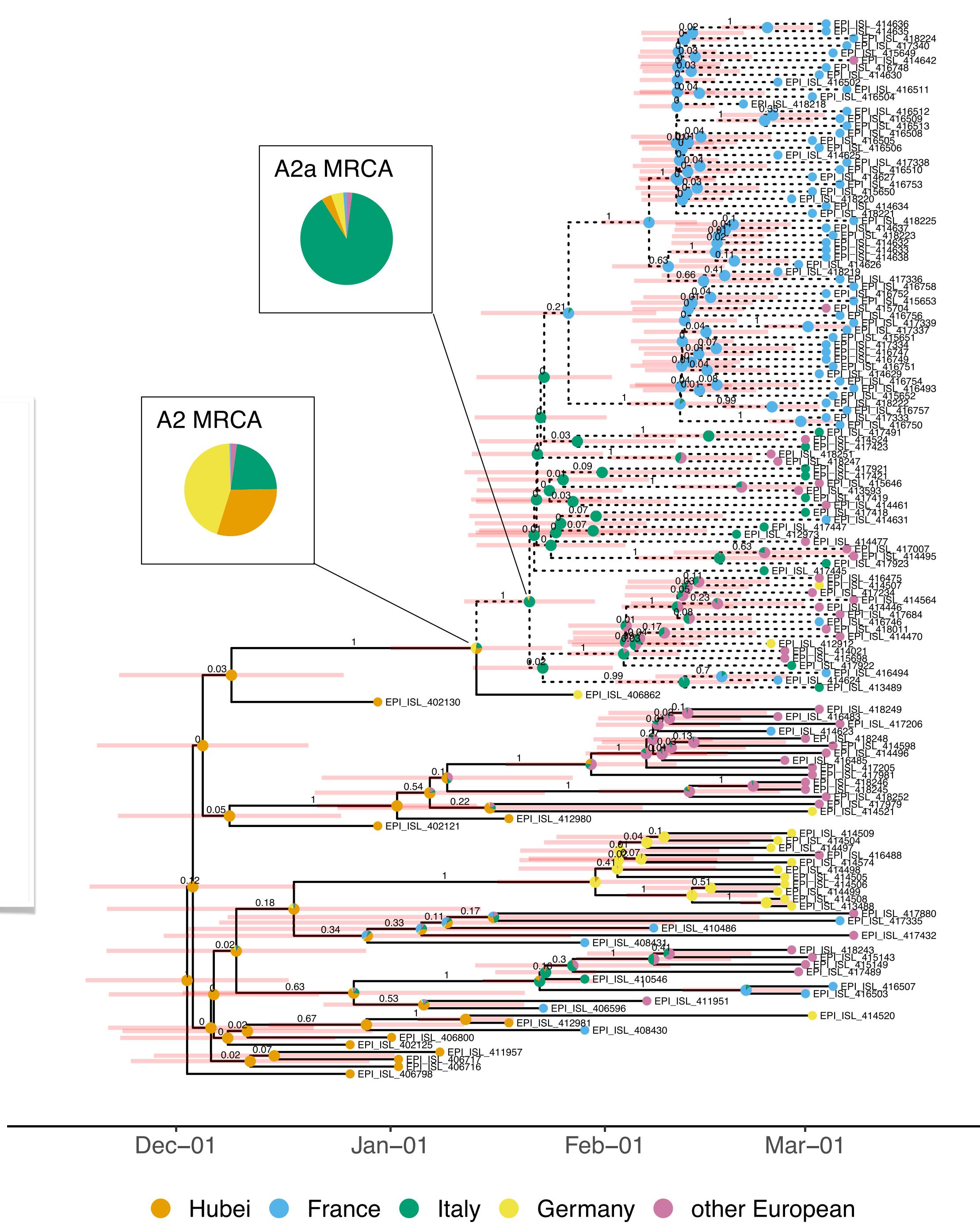
Phylogenies have been used to argue dinosaurs were incline prior to the KPg

FBD analyses suggest that we can not currently answer that question using phylogenies

# Models that include migration



The origin and early spread of SARS-CoV-2 in Europe  
Nadeau et al. 2021. PNAS



# Bayesian divergence time estimation

## The data

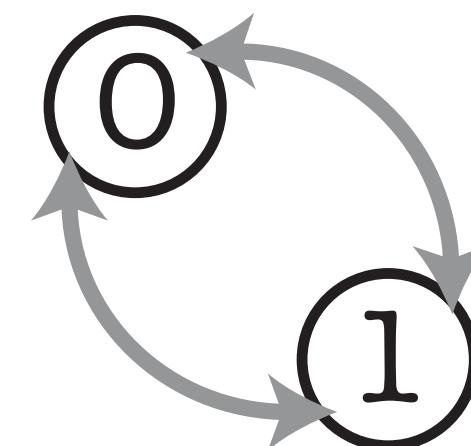
and / or  
0101... ATTG...  
1101... TTGC...  
0100... ATTC...



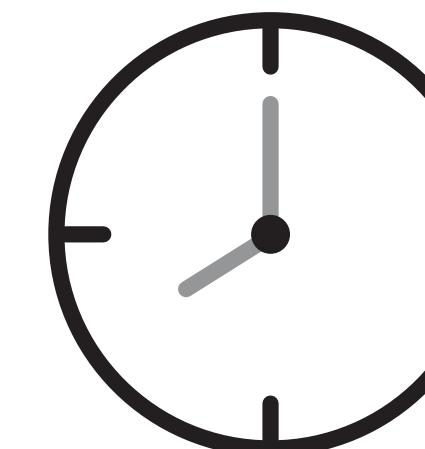
phylogenetics  
characters

sample  
ages

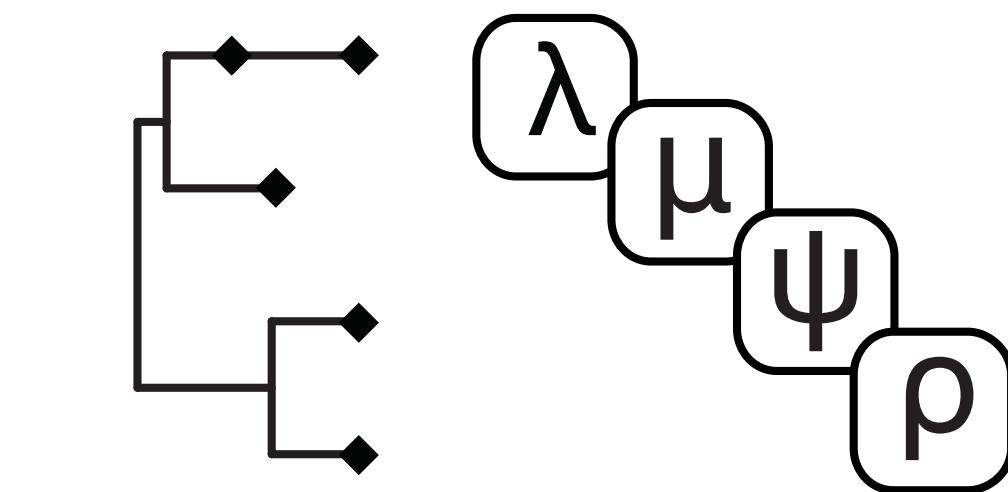
## 3 model components



substitution  
model



clock  
model



tree and tree  
model

# Using PCMs for dating

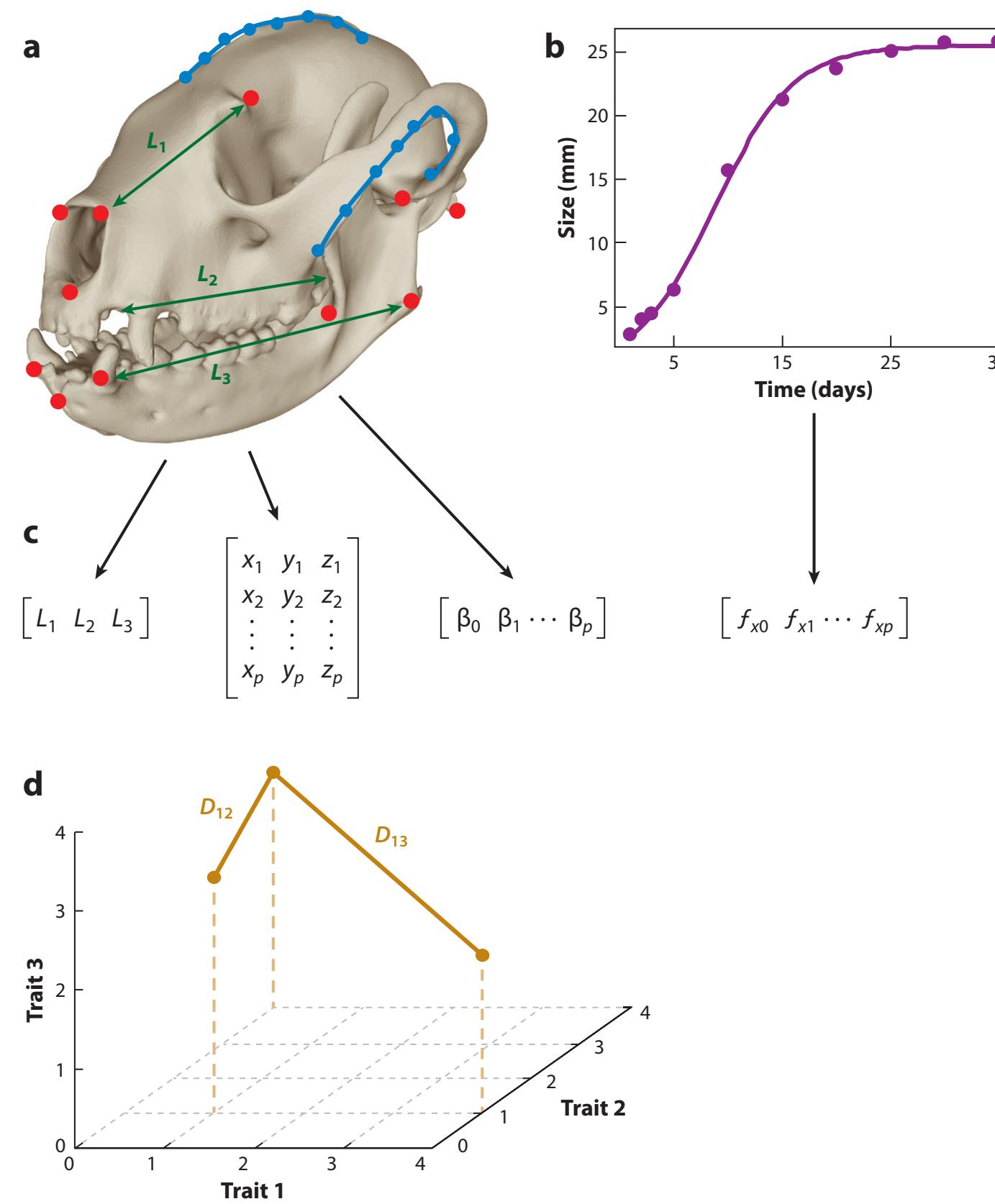
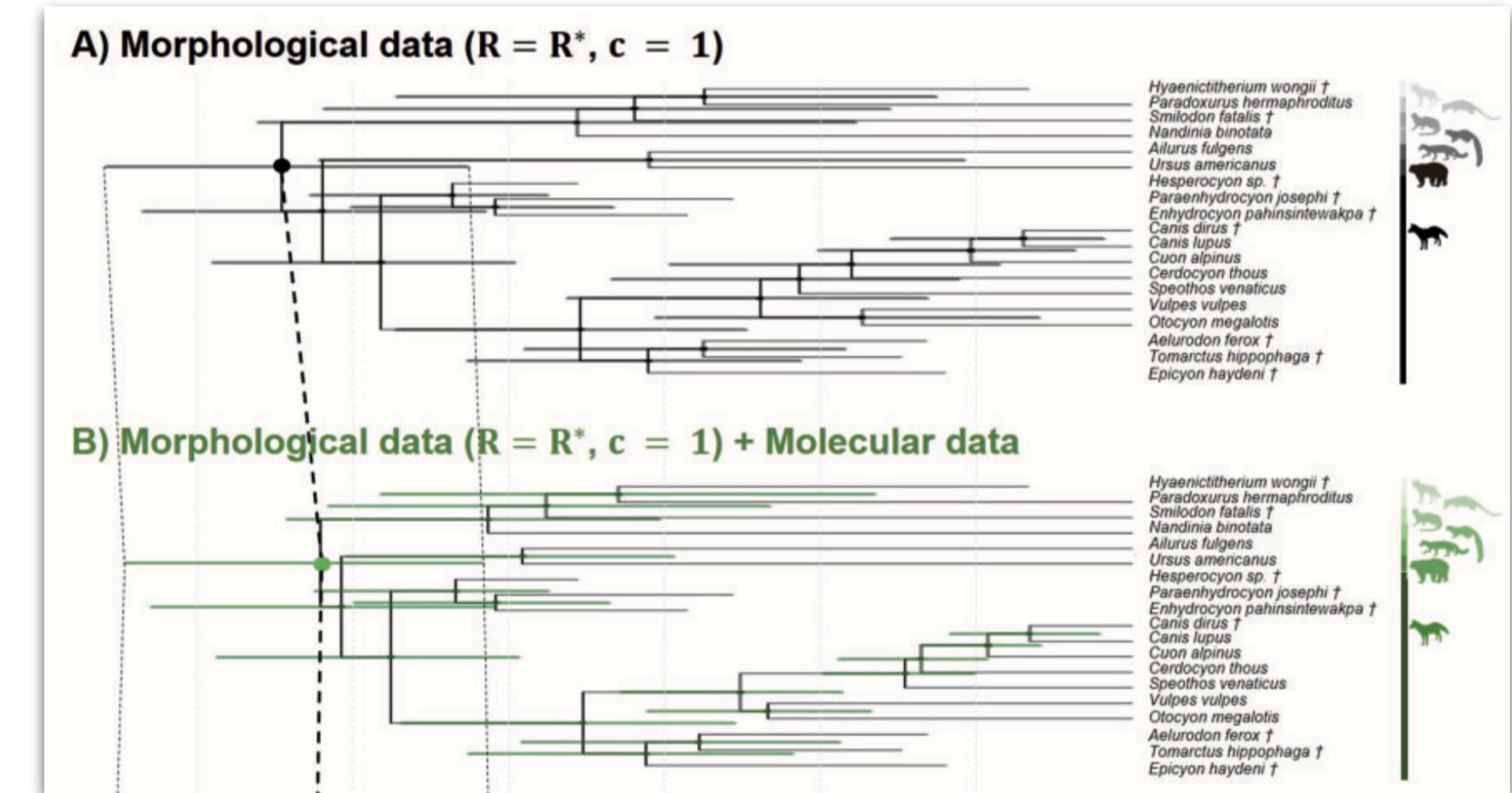
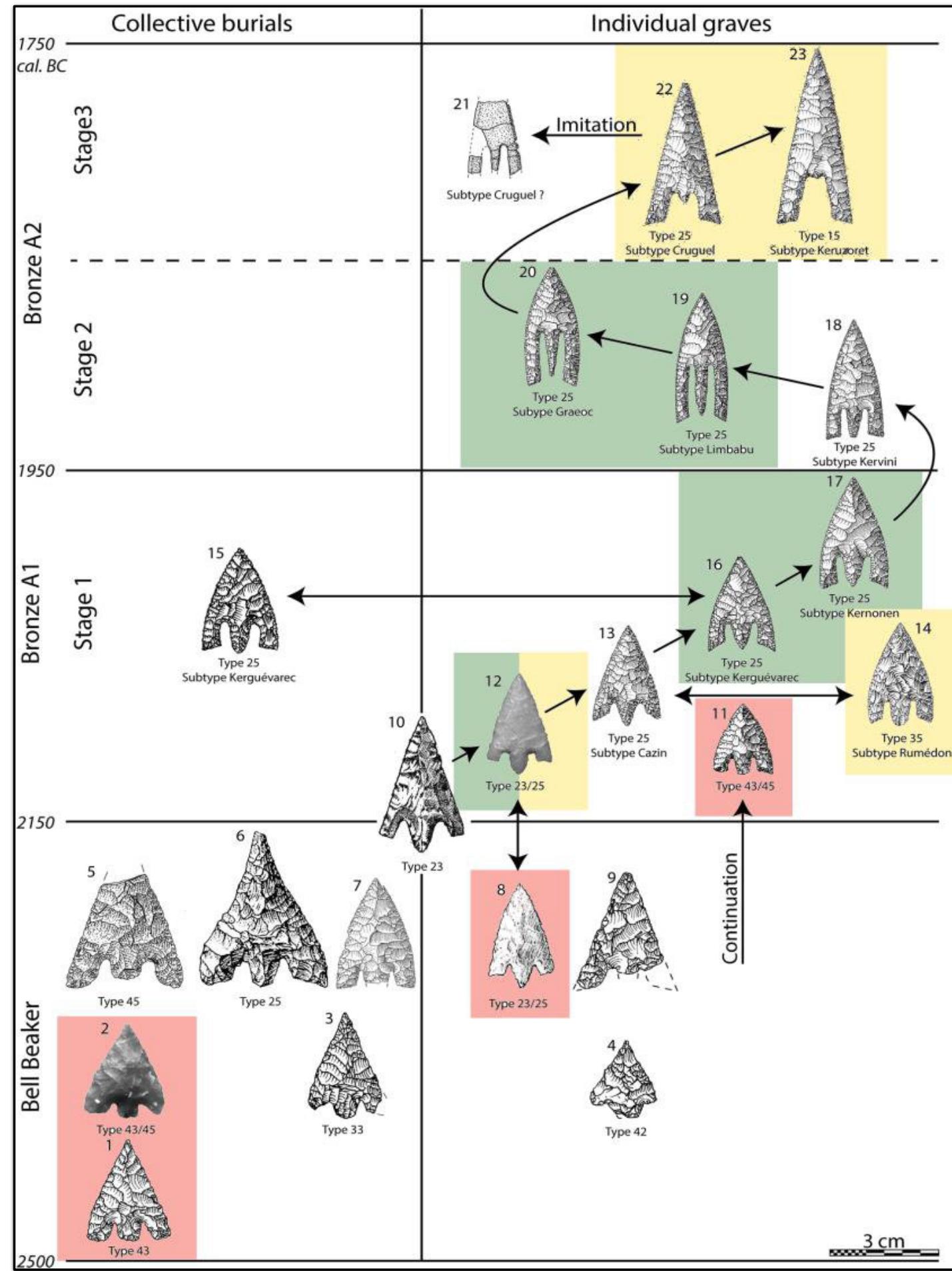


Image source Adams & Collyer ([2019](#))



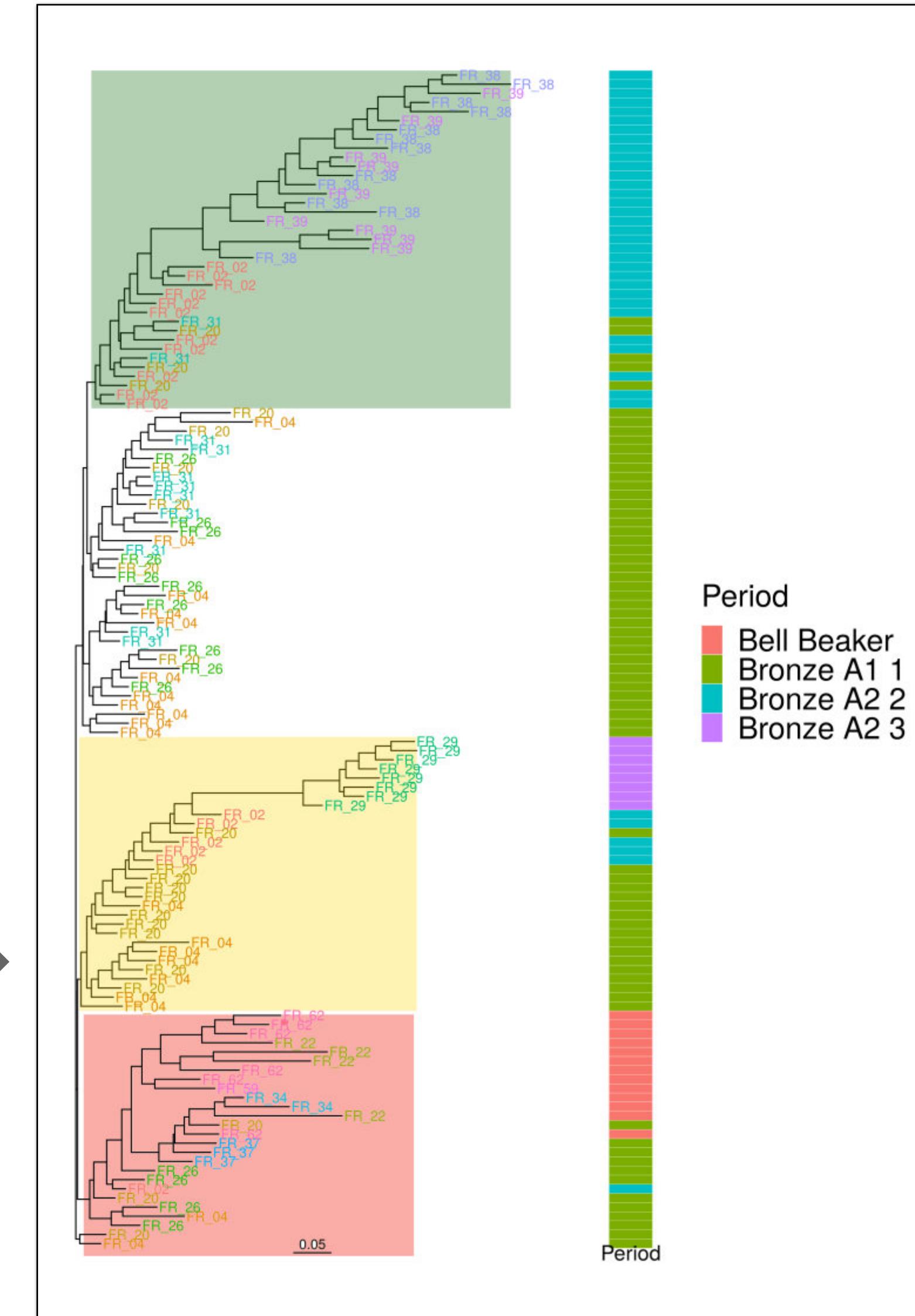
Álvarez-Carretero et al. ([2019](#)) Bayesian Estimation of Species Divergence Times Using Correlated Quantitative Characters

# Cultural evolution



← Typo-Chronology of  
Palaeolithic stone tools

Outline based NJ tree →

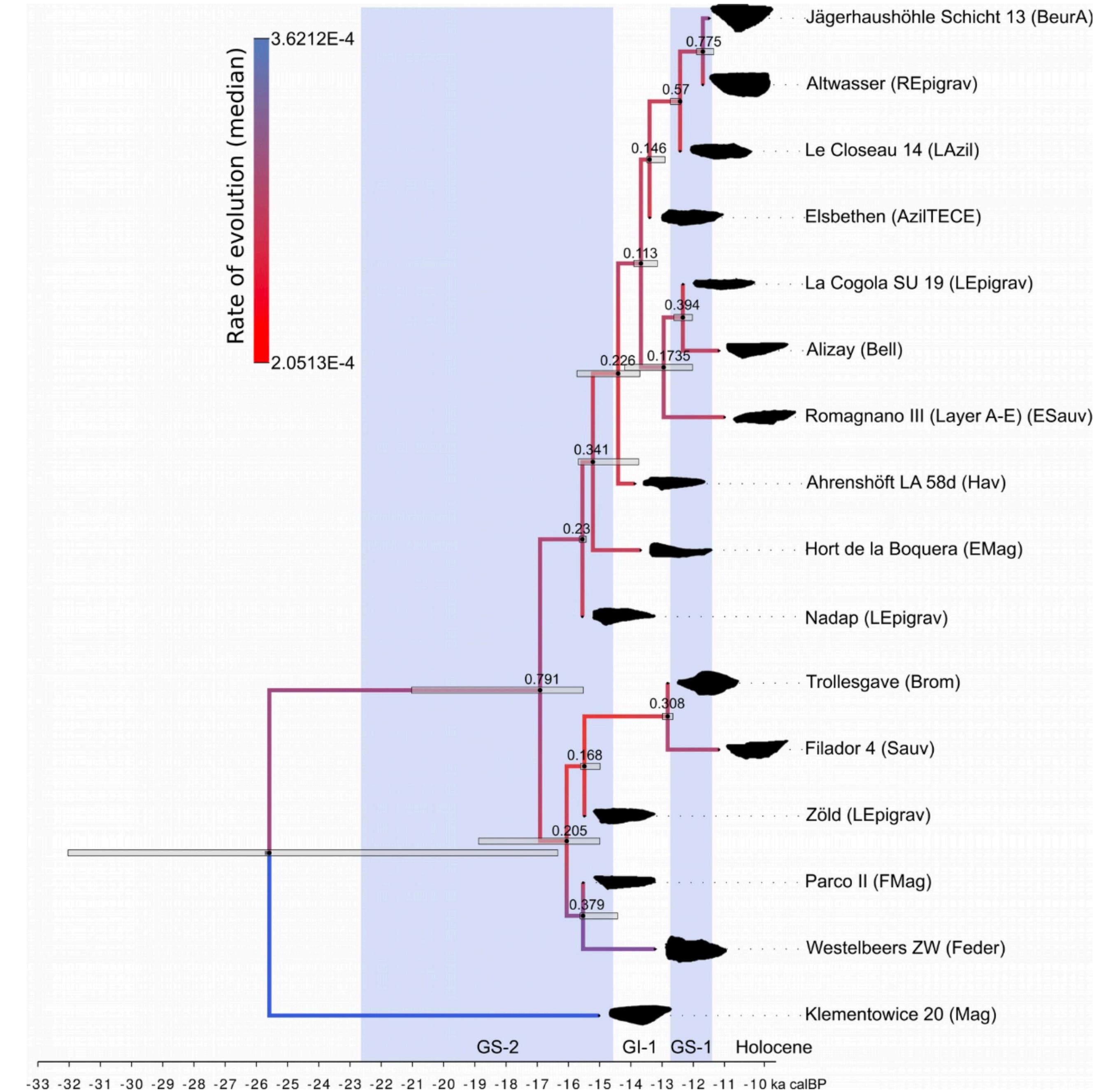


Matzig et al. 2021.

After Nicolas (2017)

# The tree topology of stone tools exhibits a lot of uncertainty

Matzig et al. (in review) A macroevolutionary analysis of European Late Upper Palaeolithic stone tool shape using a Bayesian phylodynamic framework (preprint available)



# Decoding Genomes: From Sequences to Phyldynamics

*Tanja Stadler, Carsten Magnus,  
Timothy Vaughan, Joëlle Barido-Sottani,  
Veronika Bošková, Jana S. Huisman,  
Jūlija Pečerska*

**Illustrated by Cecilia Valenzuela Agúí**

**Edited by Jūlija Pečerska**

## Obtaining the book

You will shortly be able to purchase a hard copy from Amazon. (Quality testing currently in progress.)

Alternatively, you can [Download](#) the complete PDF of the book free of charge. (See below for license information.)

## About the book

**Decoding Genomes** demonstrates how to uncover information about past evolutionary and population dynamic processes based on genomic samples. The last decades have seen considerable theoretical and methodological advances in this area. These enable the assessment of critical scientific questions such as the impact of environmental changes on biodiversity and the evolution of pathogens during recent epidemics. The book gives the reader a detailed understanding of the whole process: from genome sampling to obtaining biological insights by applying sophisticated statistical and computational analyses. In particular, sequencing of genomic samples, the alignment of sequences, molecular evolution models, phylogenetics, and phyldynamics are core topics. Statistical and computational approaches discussed include dynamic programming, maximum likelihood, Bayesian statistics, and model selection, to name a few. The concepts introduced and applied throughout the book enable readers to answer

