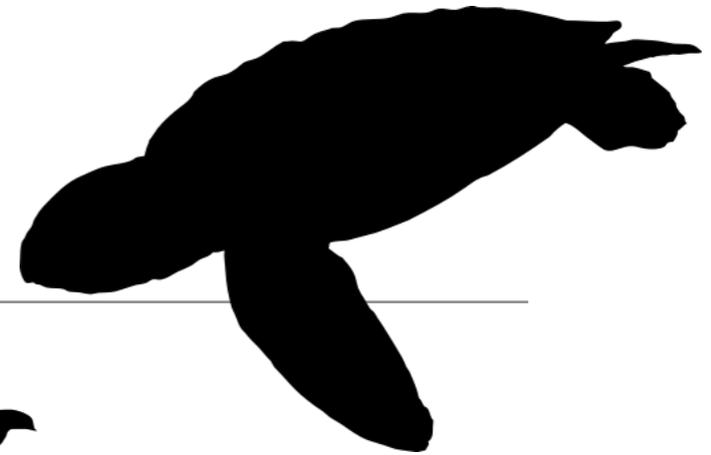


Divergence time estimation

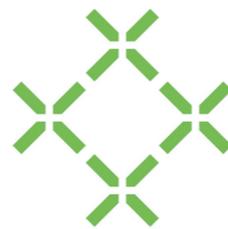
Rachel Warnock



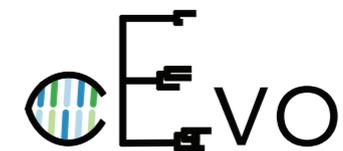
001011

ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



D-BSSE



Computational Evolution
<http://www.bsse.ethz.ch/cevo>

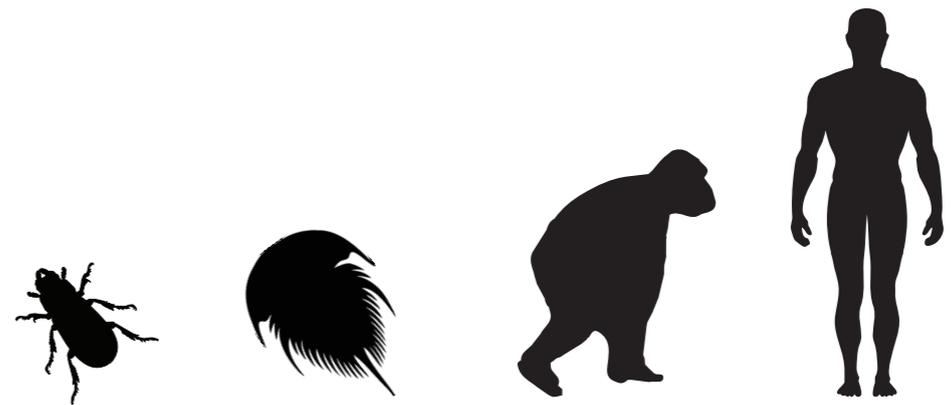
Testing hypotheses in evolutionary biology & macroevolution

t speciation or extinction times

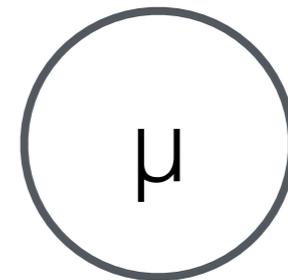
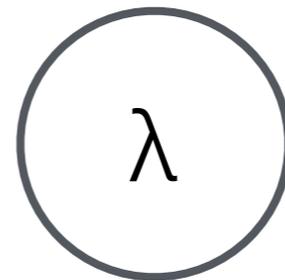
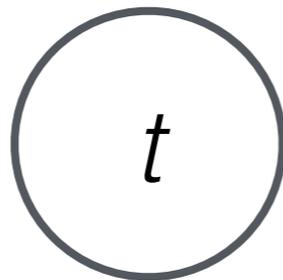
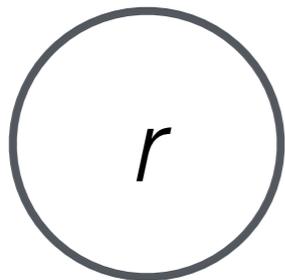
r rates of morphological or molecular evolution

λ rate of speciation

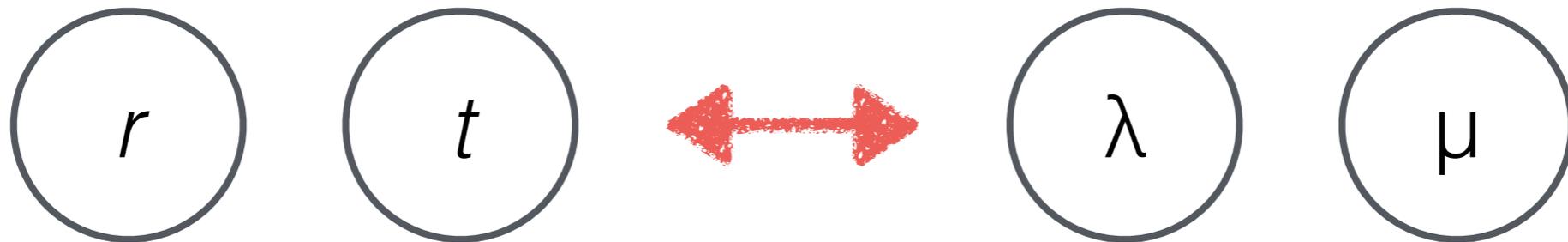
μ rate of extinction



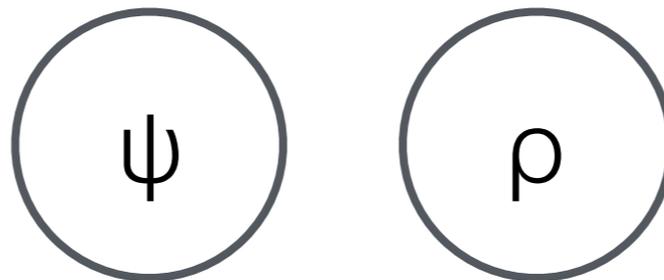
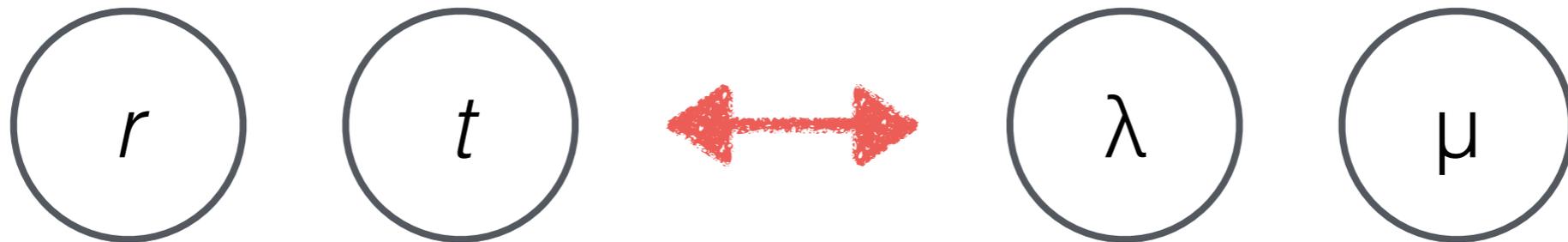
Macroevolutionary parameters of interest are
phylogenetic parameters



Macroevolutionary parameters of interest are phylogenetic parameters



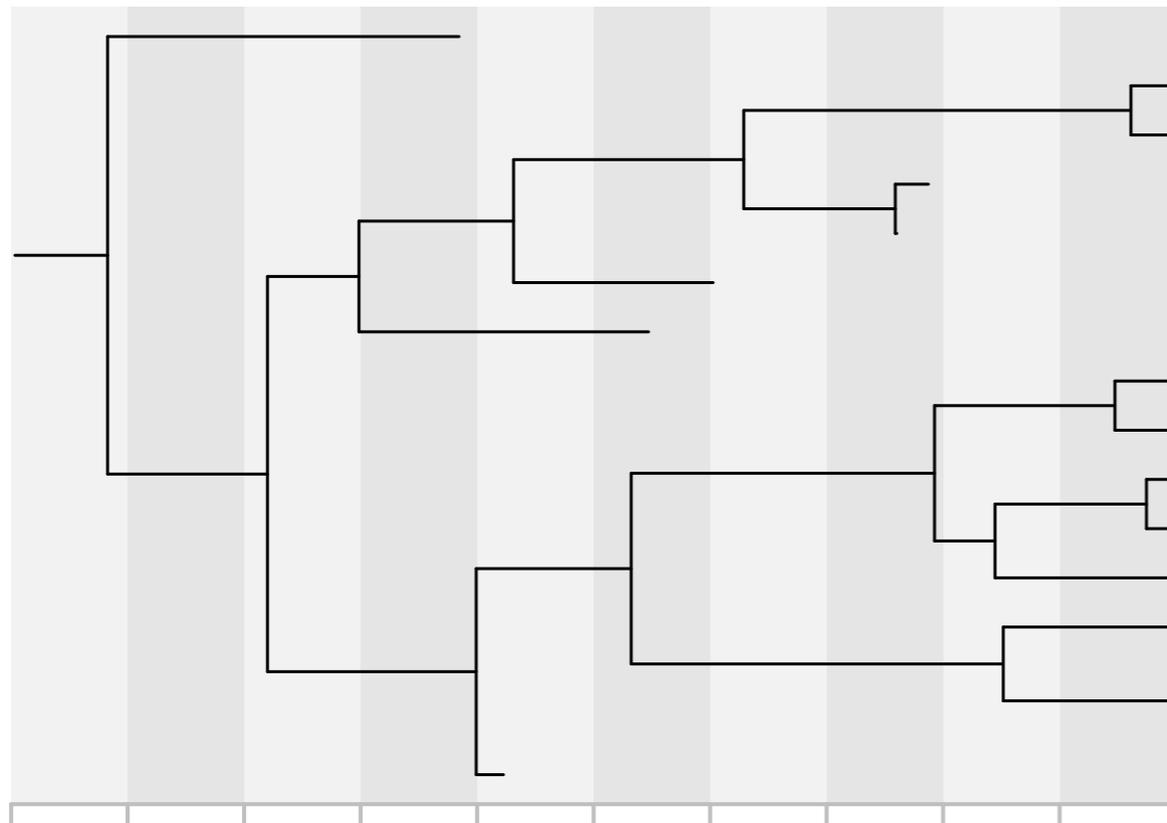
Macroevolutionary parameters of interest are phylogenetic parameters



fossil sampling rate

extant species sampling

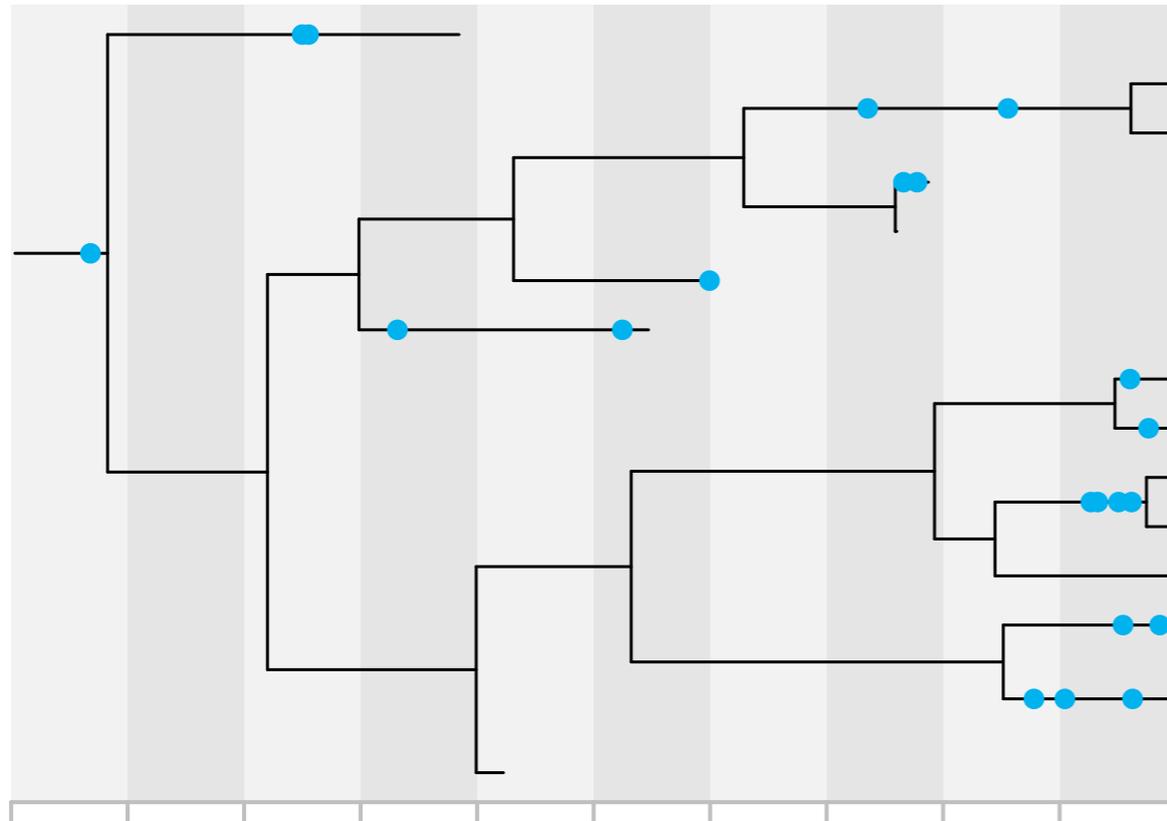
A straightforward model of evolution and sampling



λ — speciation rate

μ — extinction rate

A straightforward model of evolution and sampling

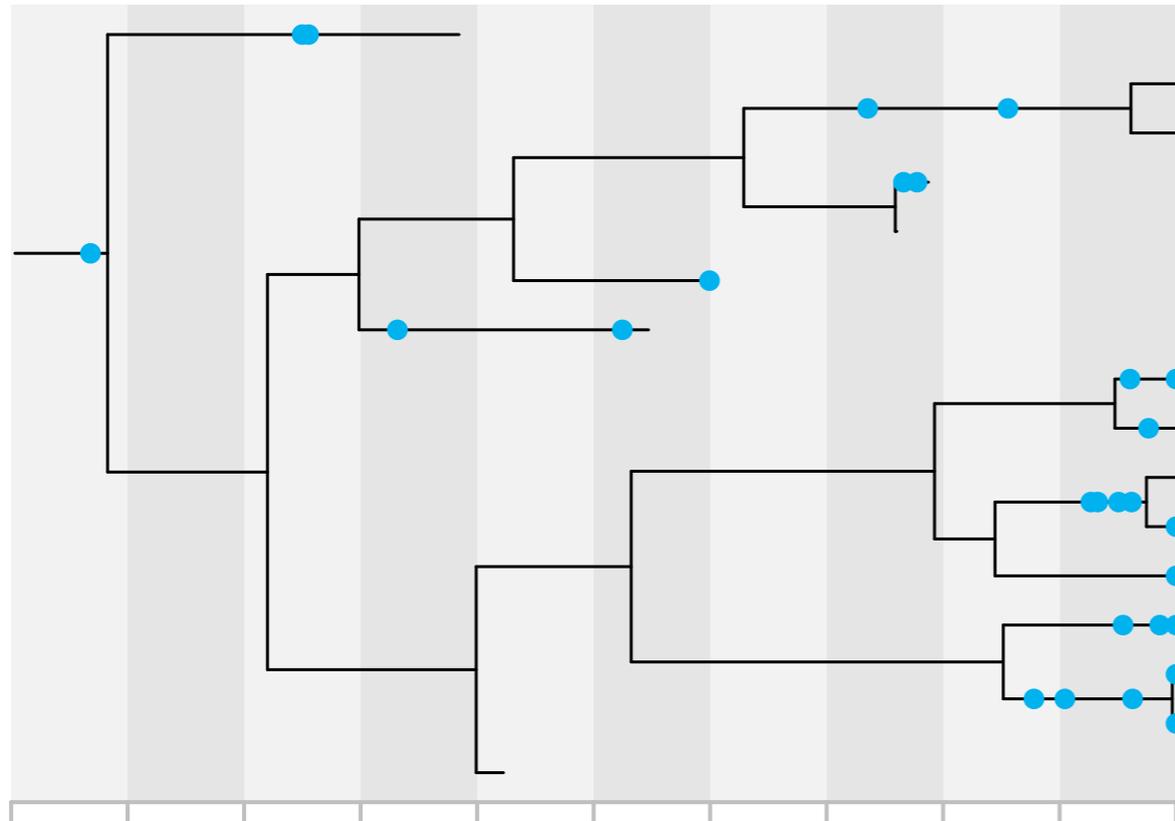


λ — speciation rate

μ — extinction rate

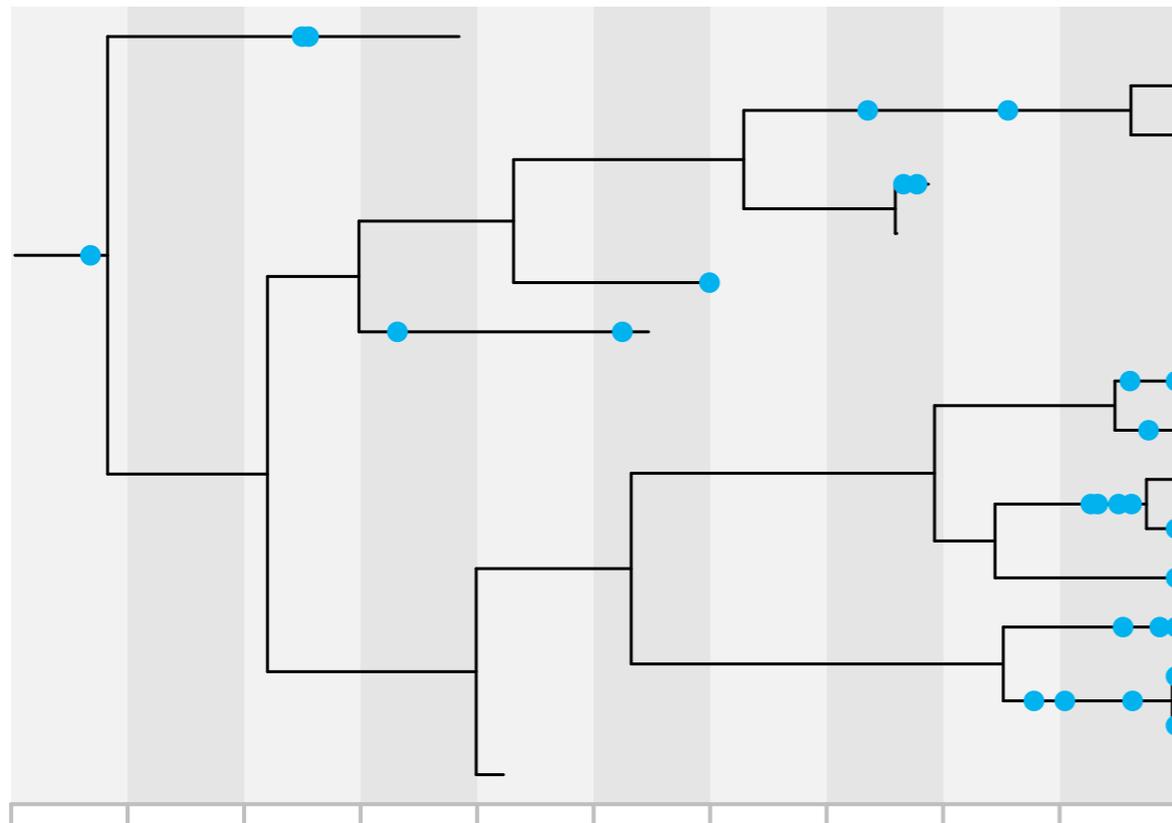
ψ — fossil sampling rate

A straightforward model of evolution and sampling



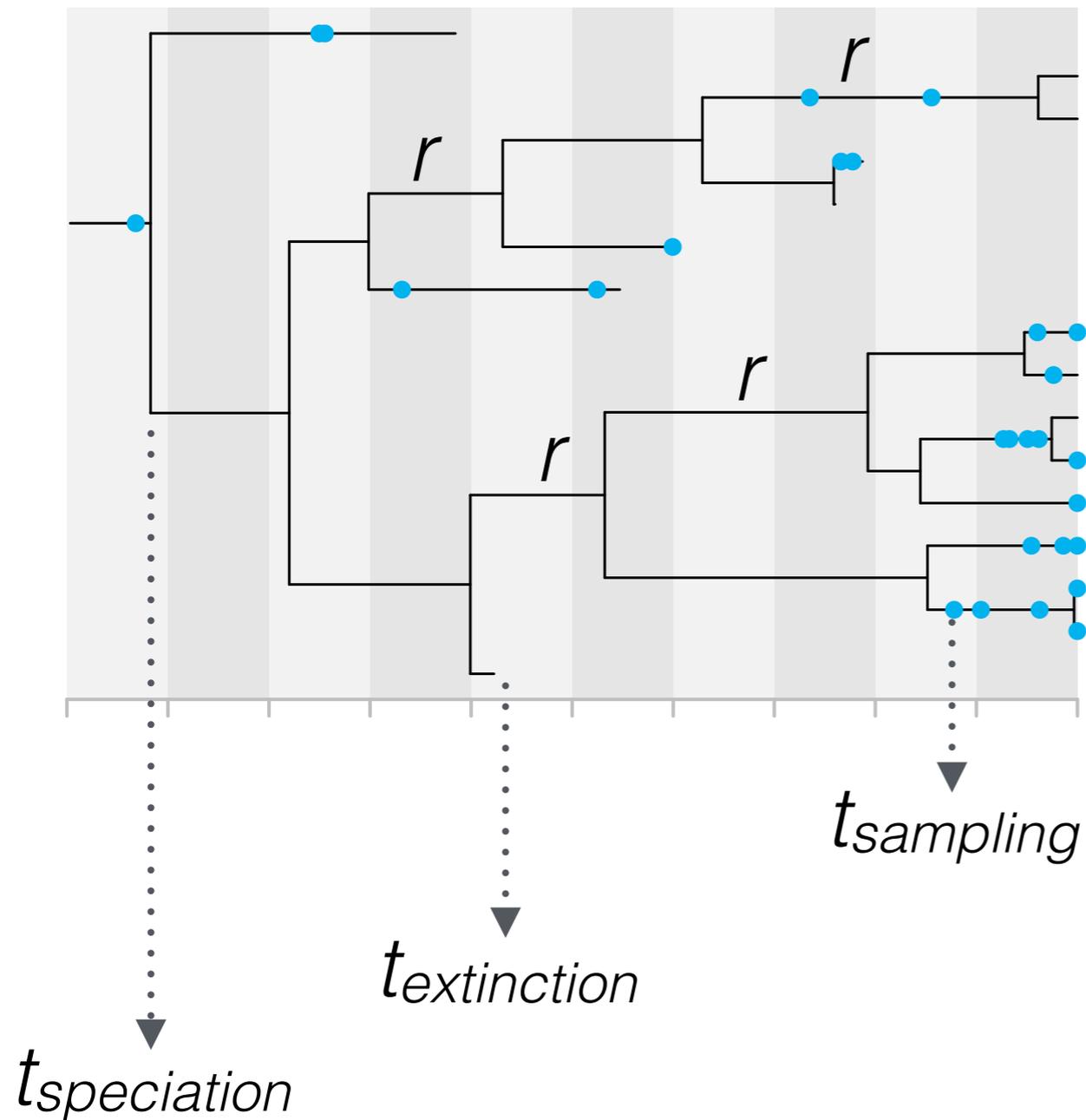
- λ — speciation rate
- μ — extinction rate
- ψ — fossil sampling rate
- ρ — extant species sampling

The fossilized birth-death process



- λ — speciation rate
- μ — extinction rate
- ψ — fossil sampling rate
- ρ — extant species sampling

A straightforward model of evolution and sampling



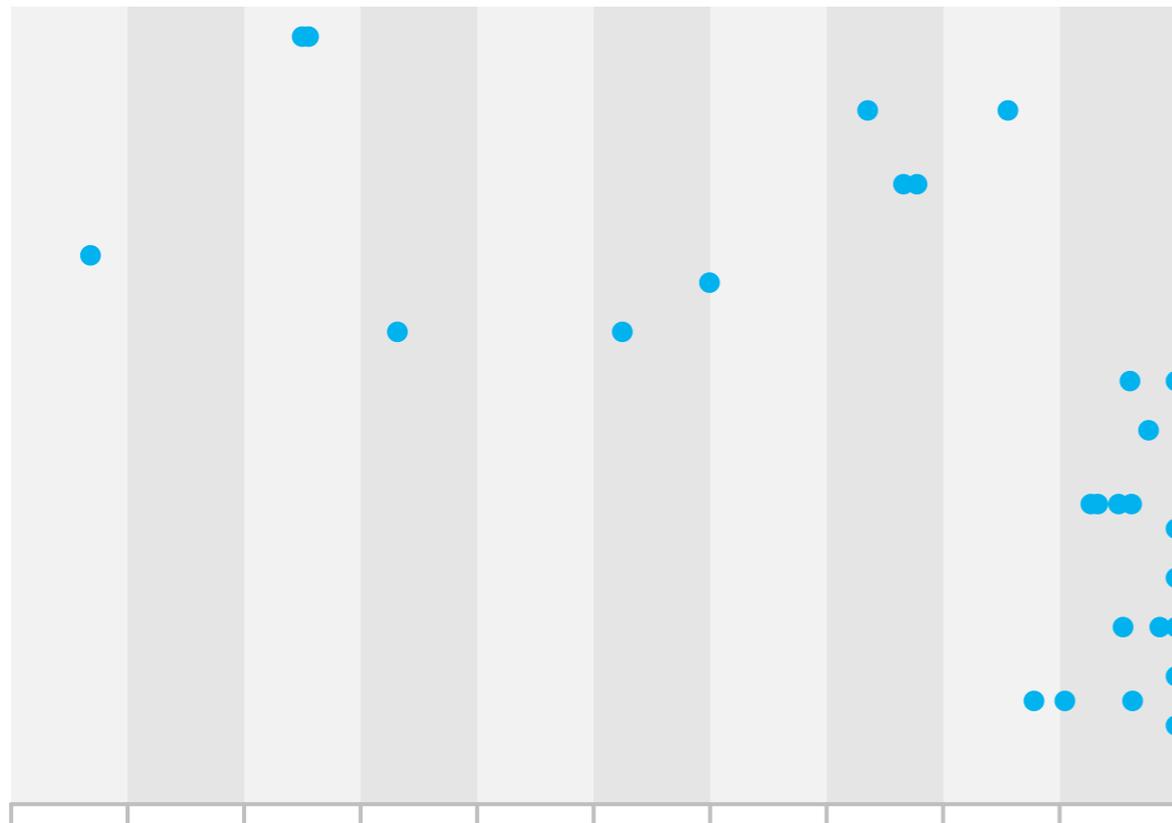
λ — speciation rate

μ — extinction rate

ψ — fossil sampling rate

ρ — extant species sampling

A straightforward model of evolution and sampling



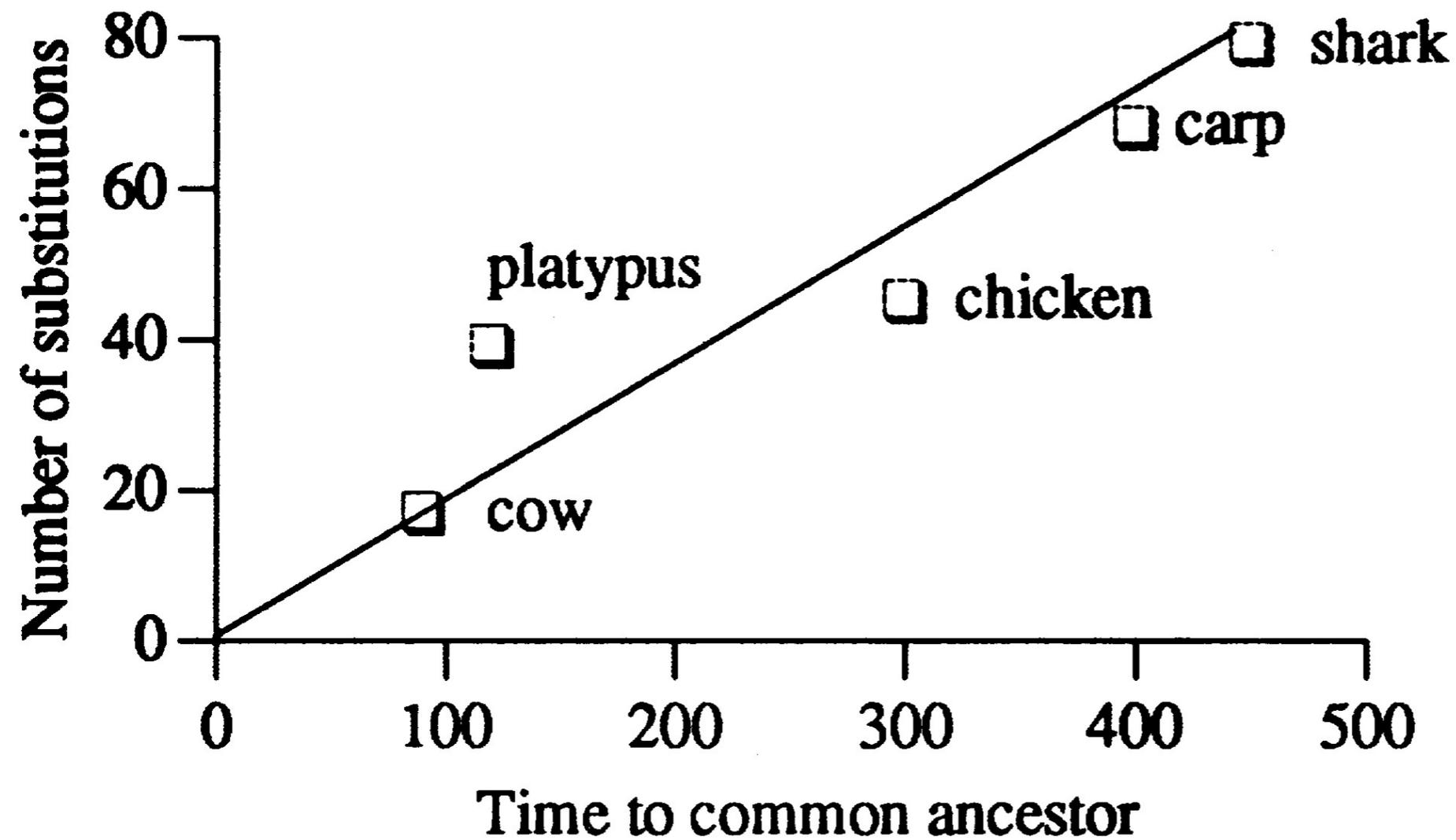
λ — speciation rate

μ — extinction rate

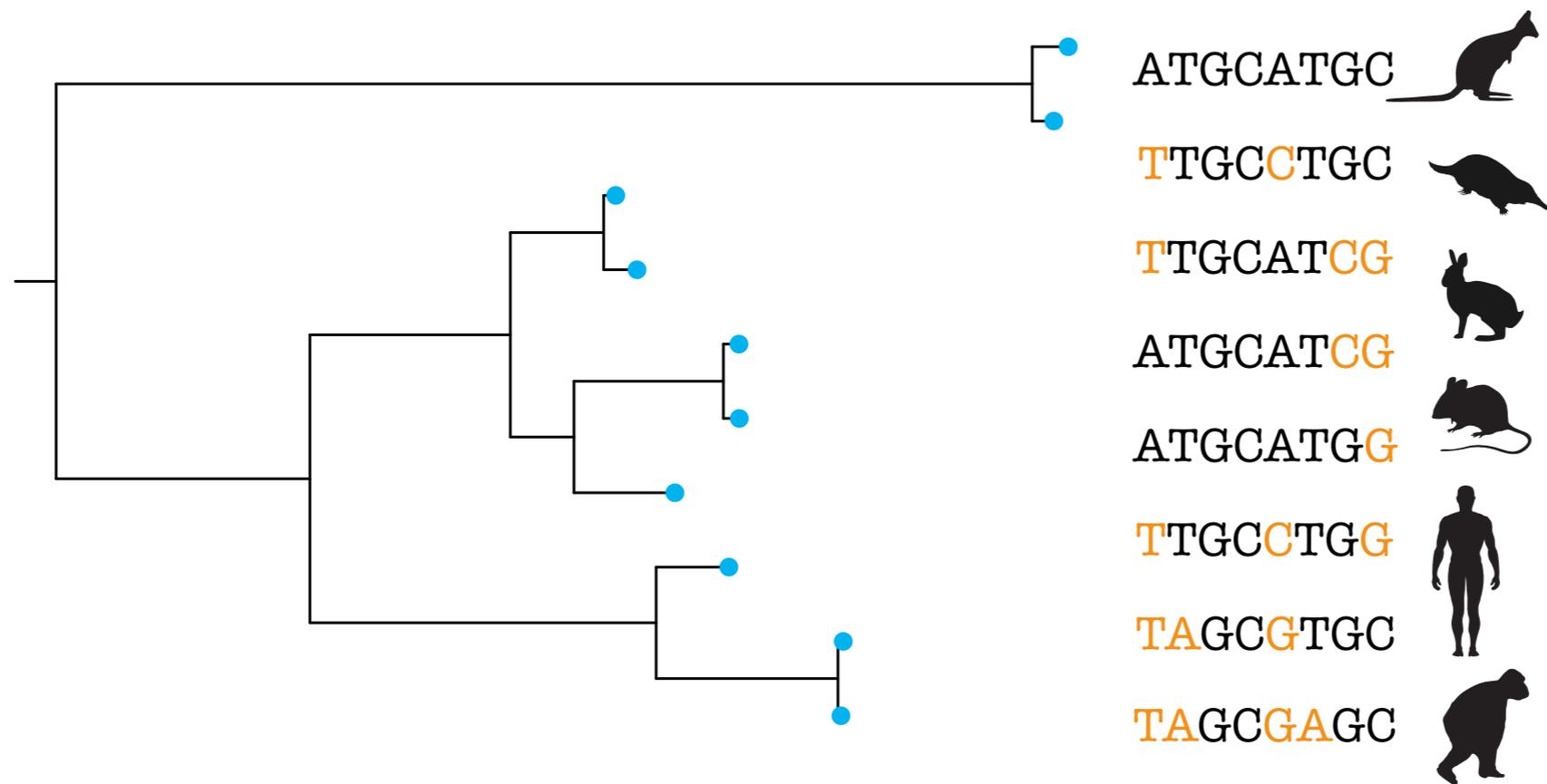
ψ — fossil sampling rate

ρ — extant species sampling

The molecular clock hypothesis



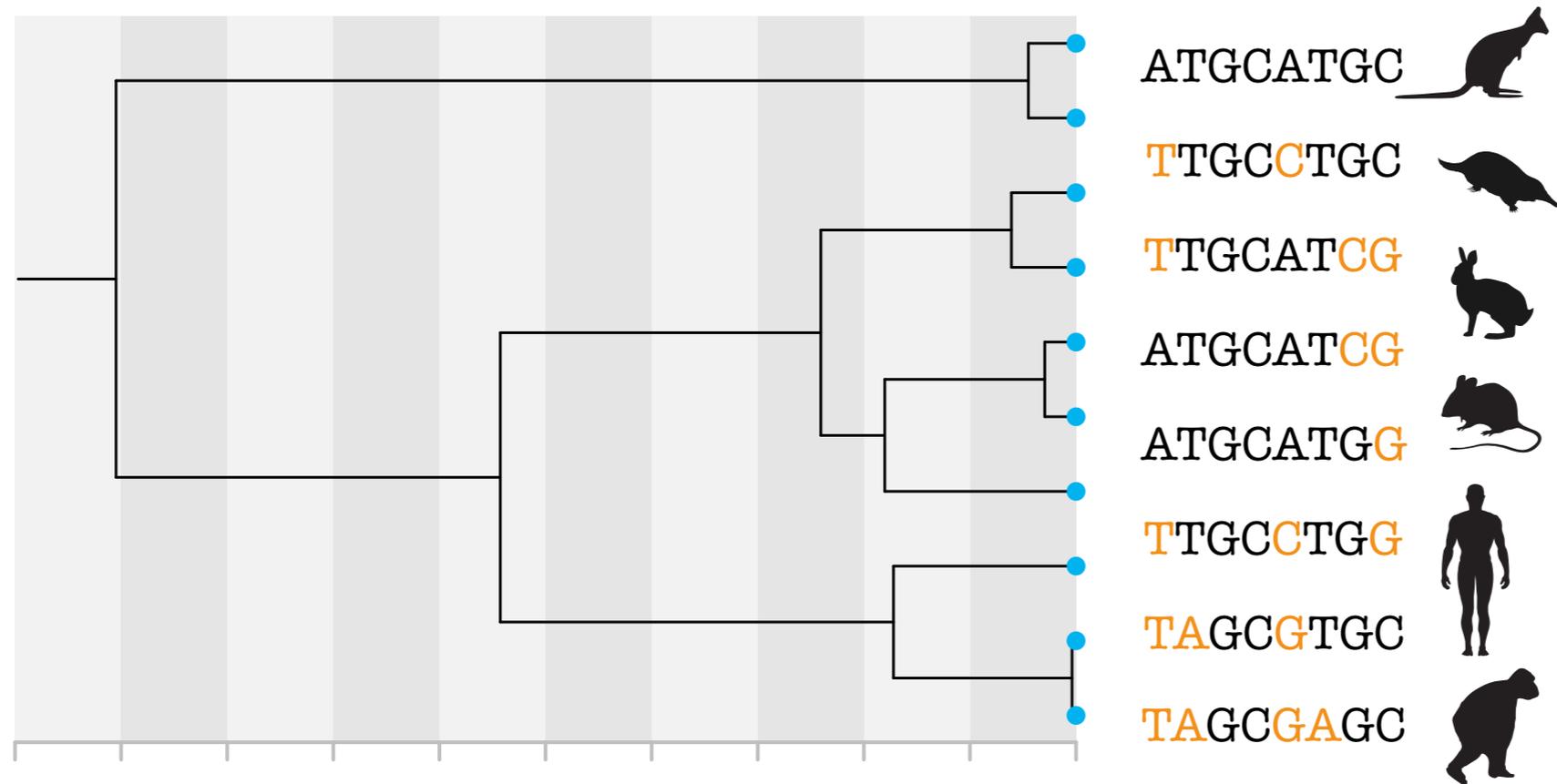
Building the tree of life



branch lengths = rate x time

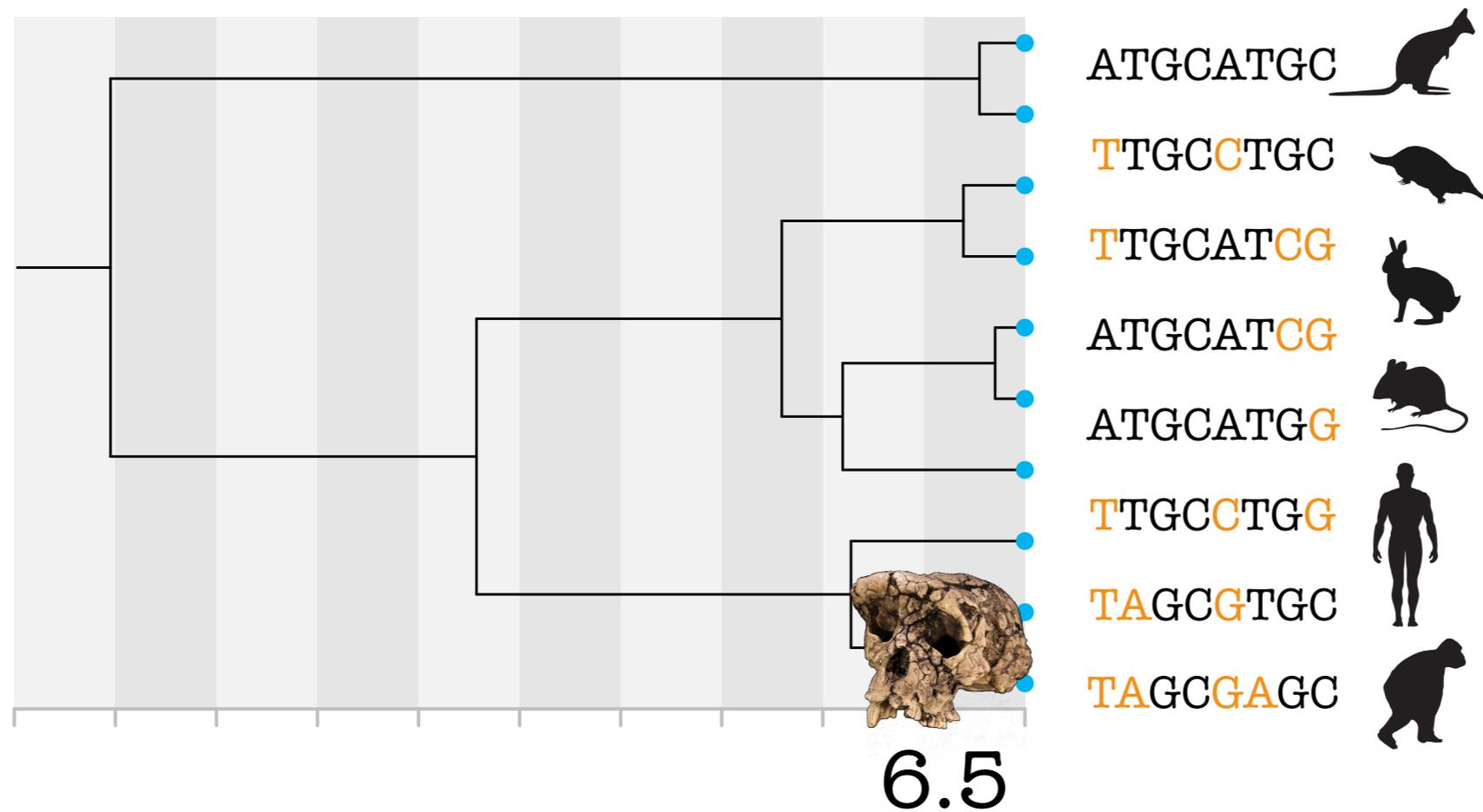
$$v = rt$$

Dating the tree of life



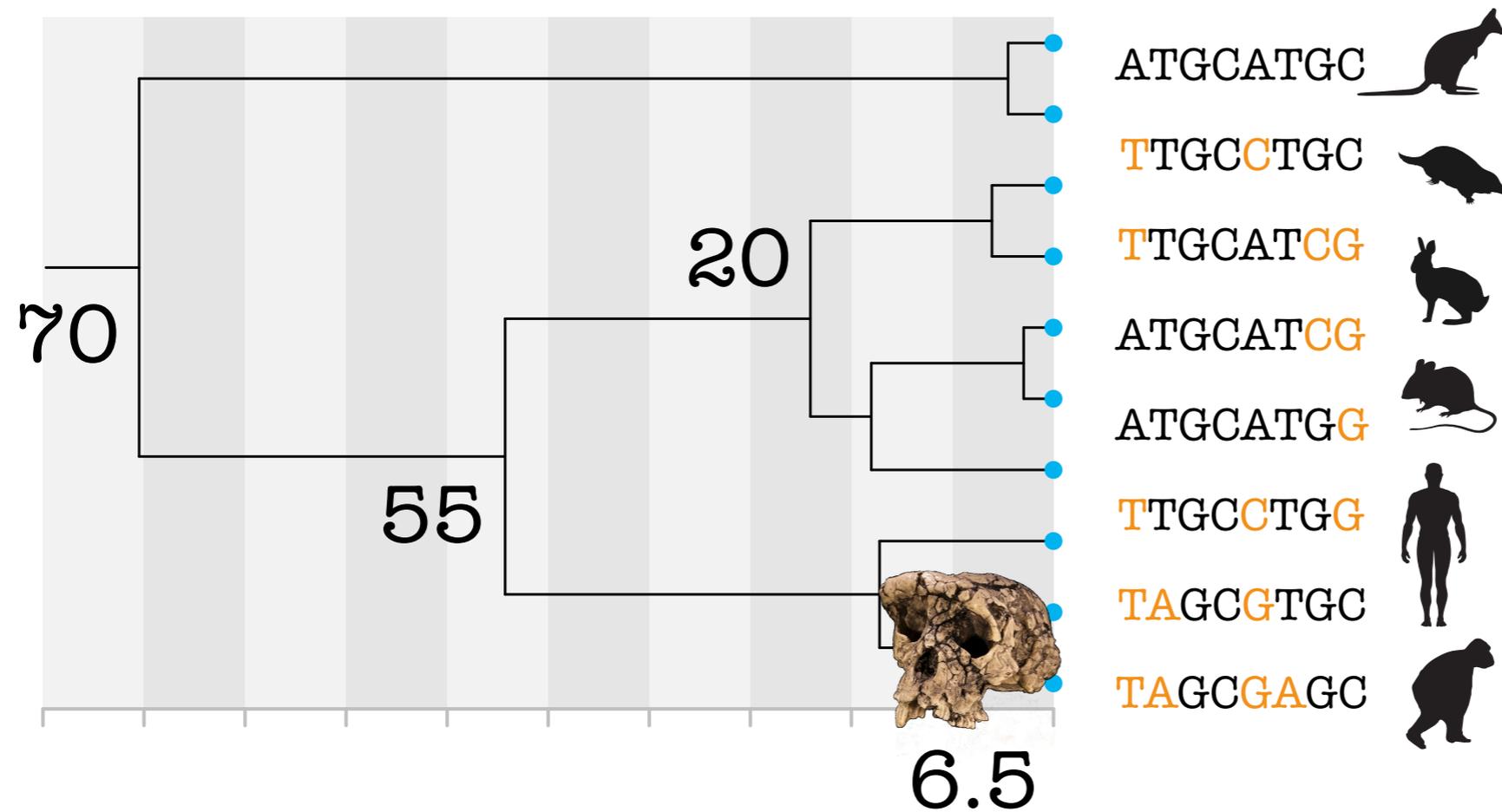
branch lengths = time

Calibrating the molecular clock



branch lengths = time

Calibrating the molecular clock



branch lengths = time

Bayesian divergence time estimation

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

Bayesian divergence time estimation

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

likelihood

priors

posterior

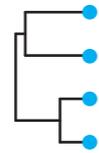
marginal likelihood of the data

The diagram illustrates the Bayesian formula for divergence time estimation. The formula is $P(\text{model} \mid \text{data}) = \frac{P(\text{data} \mid \text{model}) P(\text{model})}{P(\text{data})}$. The terms are color-coded: 'model' is blue, 'data' is orange, and 'likelihood' is black. Red arrows point from the labels to the corresponding terms in the formula: 'likelihood' points to $P(\text{data} \mid \text{model})$, 'priors' points to $P(\text{model})$, 'posterior' points to $P(\text{model} \mid \text{data})$, and 'marginal likelihood of the data' points to $P(\text{data})$.

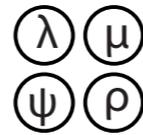
Bayesian divergence time estimation

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

alignment



tree



tree
model



calibration
priors



substitution
model



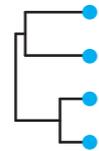
clock
model

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

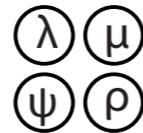
Bayesian divergence time estimation

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

alignment



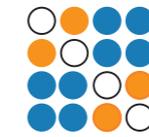
tree



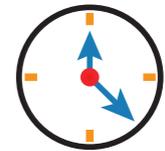
tree
model



calibration
priors



substitution
model

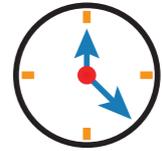
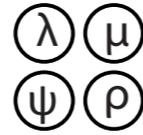
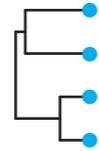


clock
model

$$P(\text{tree, model, priors, substitution model, clock model} \mid \text{ACAC... TCAC... ACAG...}) = \frac{P(\text{ACAC... TCAC... ACAG...} \mid \text{tree, model, priors, substitution model, clock model}) P(\text{tree, model, priors, substitution model, clock model})}{P(\text{ACAC... TCAC... ACAG...})}$$

Bayesian divergence time estimation

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



alignment

tree

tree
model

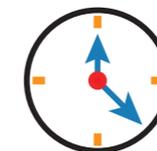
calibration
priors

substitution
model

clock
model

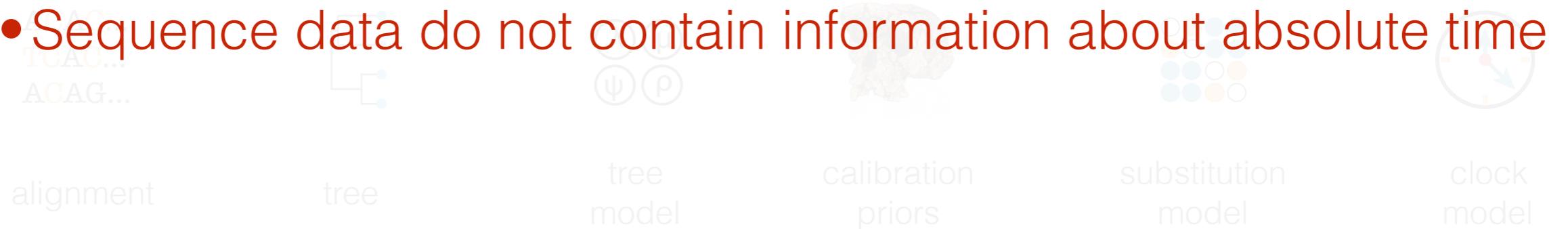
$$P\left(\begin{array}{c} \text{tree} \\ \lambda \mu \\ \psi \rho \end{array} \begin{array}{c} \text{skull} \\ \text{grid} \\ \text{clock} \end{array} \mid \begin{array}{c} \text{ACAC...} \\ \text{TCAC...} \\ \text{ACAG...} \end{array}\right) = \frac{P\left(\begin{array}{c} \text{ACAC...} \\ \text{TCAC...} \\ \text{ACAG...} \end{array} \mid \begin{array}{c} \text{tree} \\ \lambda \mu \\ \psi \rho \end{array} \begin{array}{c} \text{skull} \\ \text{grid} \\ \text{clock} \end{array}\right) P\left(\begin{array}{c} \text{tree} \\ \lambda \mu \\ \psi \rho \end{array} \begin{array}{c} \text{skull} \\ \text{grid} \\ \text{clock} \end{array}\right)}{P\left(\begin{array}{c} \text{ACAC...} \\ \text{TCAC...} \\ \text{ACAG...} \end{array}\right)}$$

Coming up...



Bayesian divergence time estimation

- Sequence data do not contain information about absolute time



$$P(\text{tree, tree model, calibration priors, substitution model, clock model} \mid \begin{matrix} \text{ACAC...} \\ \text{TCAC...} \\ \text{ACAG...} \end{matrix}) = \frac{P(\begin{matrix} \text{ACAC...} \\ \text{TCAC...} \\ \text{ACAG...} \end{matrix} \mid \text{tree, tree model, calibration priors, substitution model, clock model}) P(\text{tree, tree model, calibration priors, substitution model, clock model})}{P(\begin{matrix} \text{ACAC...} \\ \text{TCAC...} \\ \text{ACAG...} \end{matrix})}$$

Coming up...



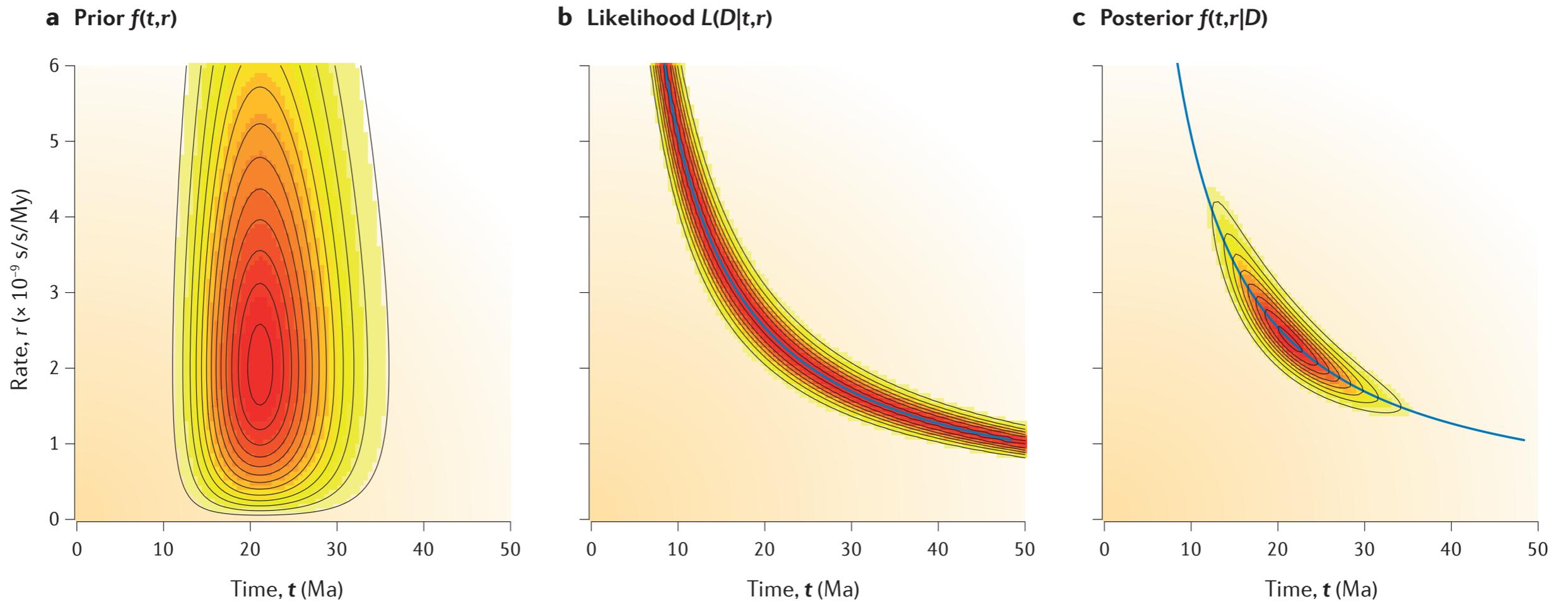
Bayesian divergence time estimation

- Sequence data do not contain information about absolute time
- This has several important consequences:
 - We need strong prior information about the divergence times or the substitution rate
 - Model selection cannot be used to select among possible calibration strategies
 - If there is uncertainty in the calibrations, even an infinite amount of sequence data won't completely eliminate uncertainty in the posteriors

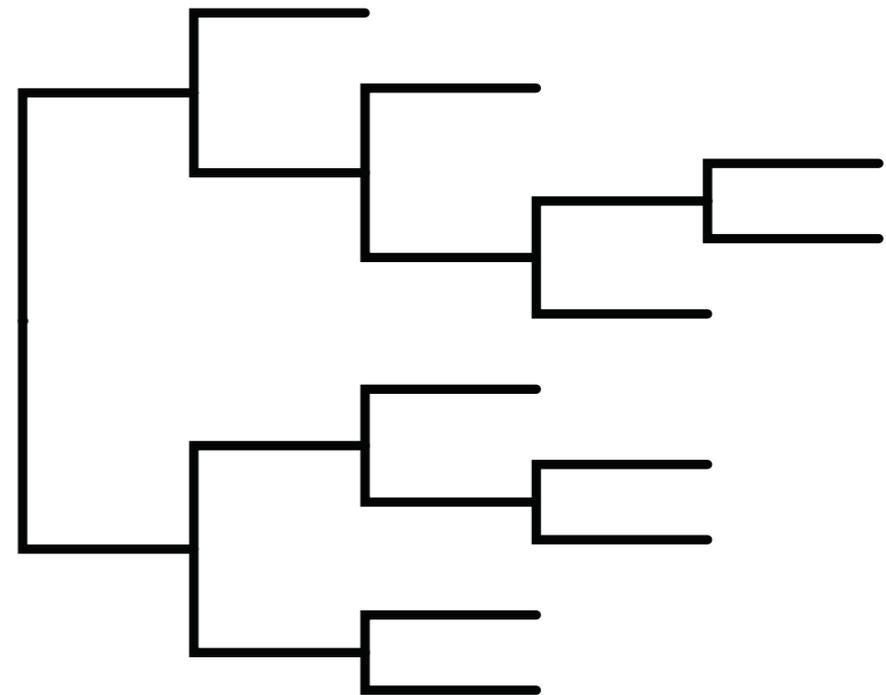
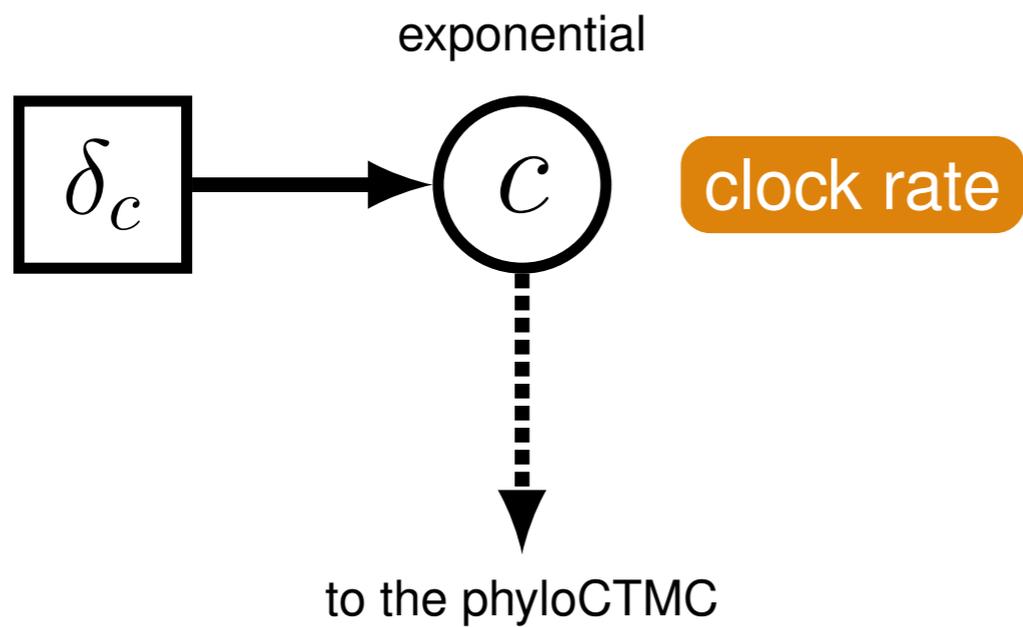
Coming up...



Rate and time are only semi-identifiable

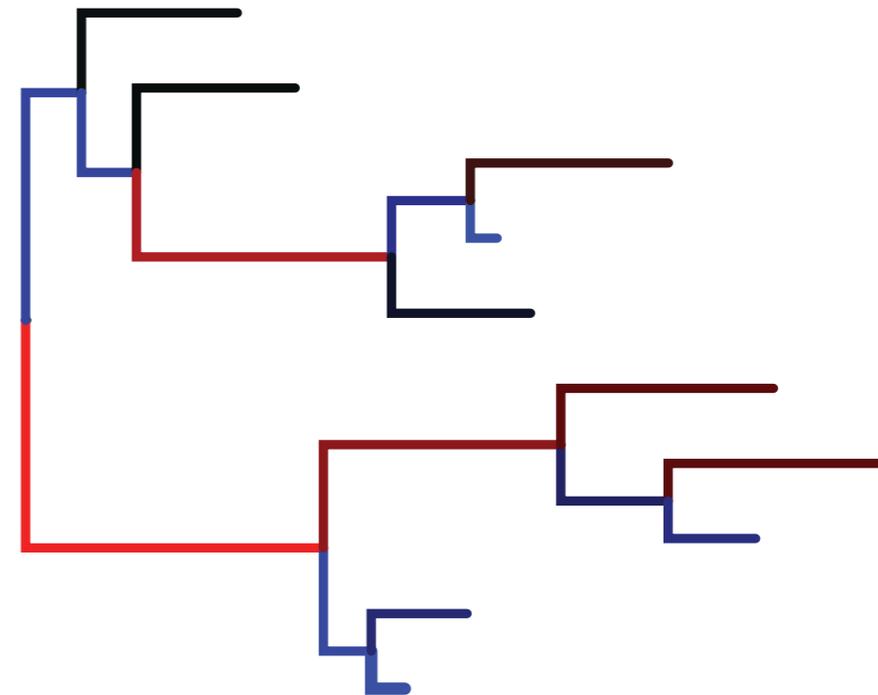
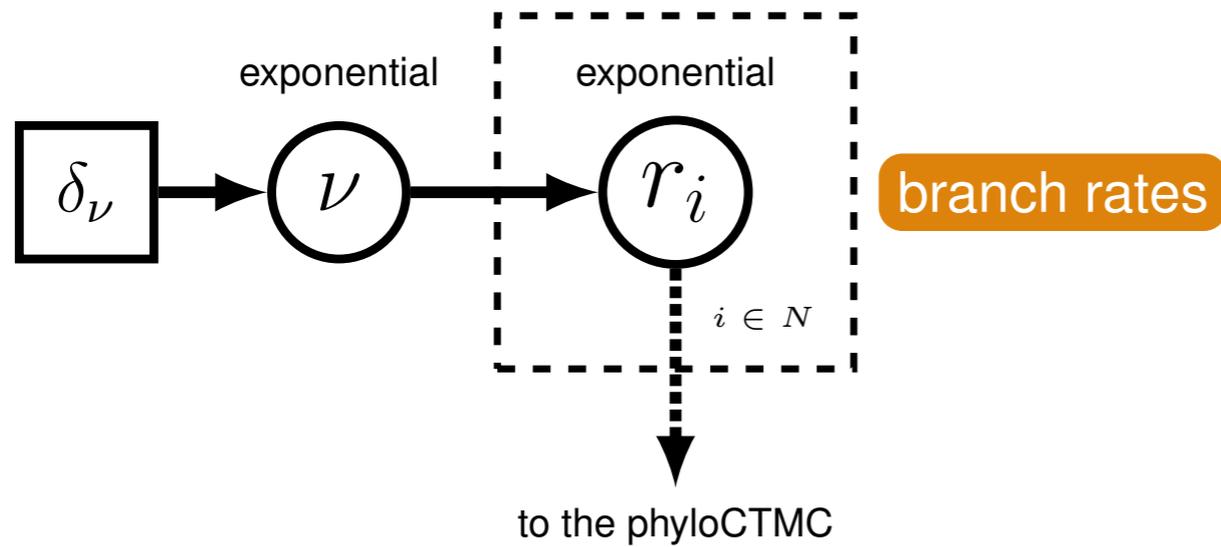


The global molecular clock model



branch length = substitution rate
low  high

The independent uncorrelated rates model



branch length = substitution rate

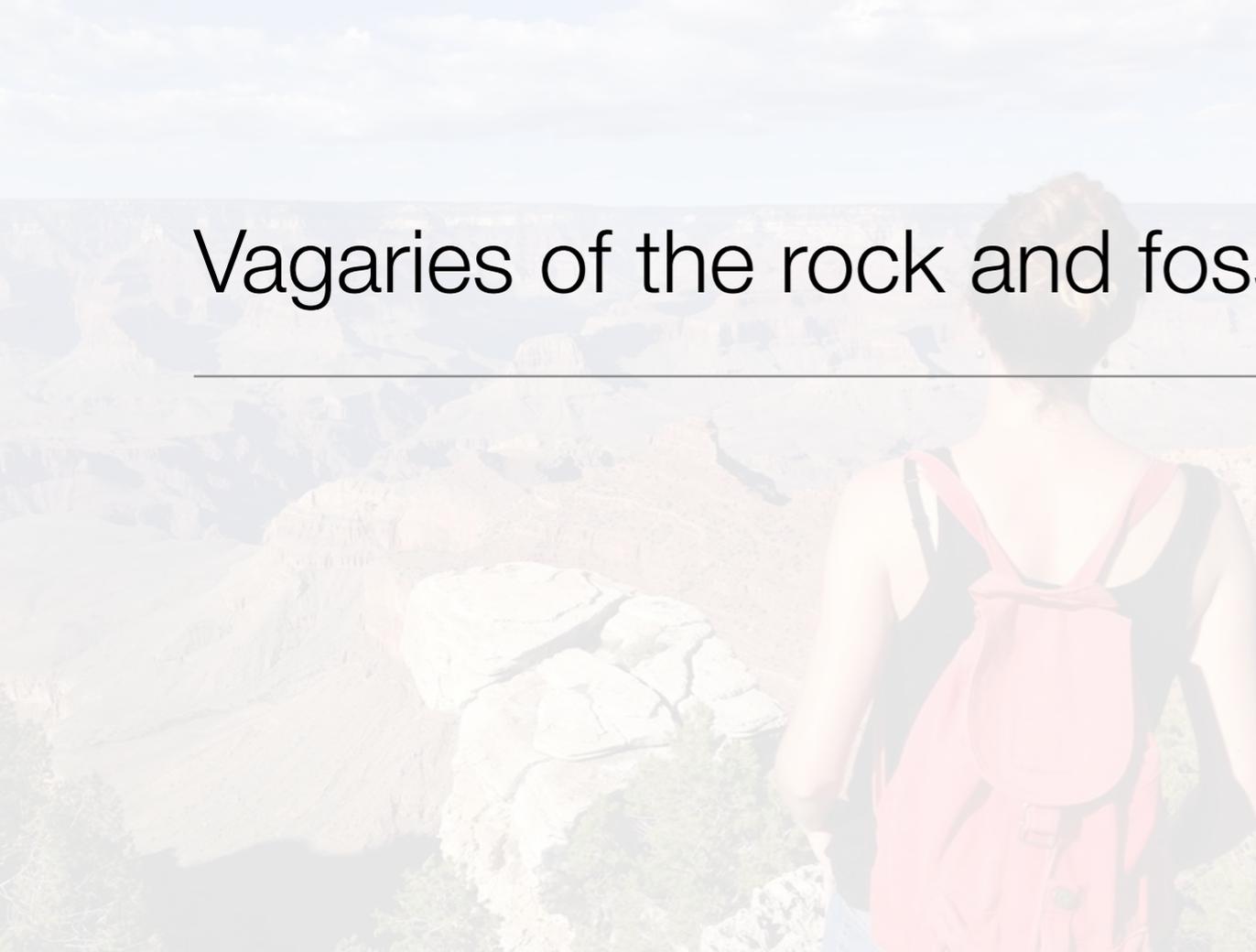
low  high

Many molecular clock models



- **Global clock** (Zuckerkandl & Pauling, 1962)
- **Uncorrelated/independent rates models** (Drummond et al. 2006; Rannala & Yang 2007; Lepage et al. 2007)
- **Log-normally distributed autocorrelated rates** (Thorne, Kishino & Painter 1998; Kishino, Thorne & Bruno 2001; Thorne & Kishino 2002)
- **Local clocks** (Hasegawa, Kishino & Yano 1989; Kishino & Hasegawa 1990; Yoder & Yang 2000; Yang & Yoder 2003, Drummond and Suchard 2010)
- **Mixture models on branch rates** (Heath, Holder, Huelsenbeck 2012)
- **Punctuated rate change model** (Huelsenbeck, Larget and Swoford 2000)

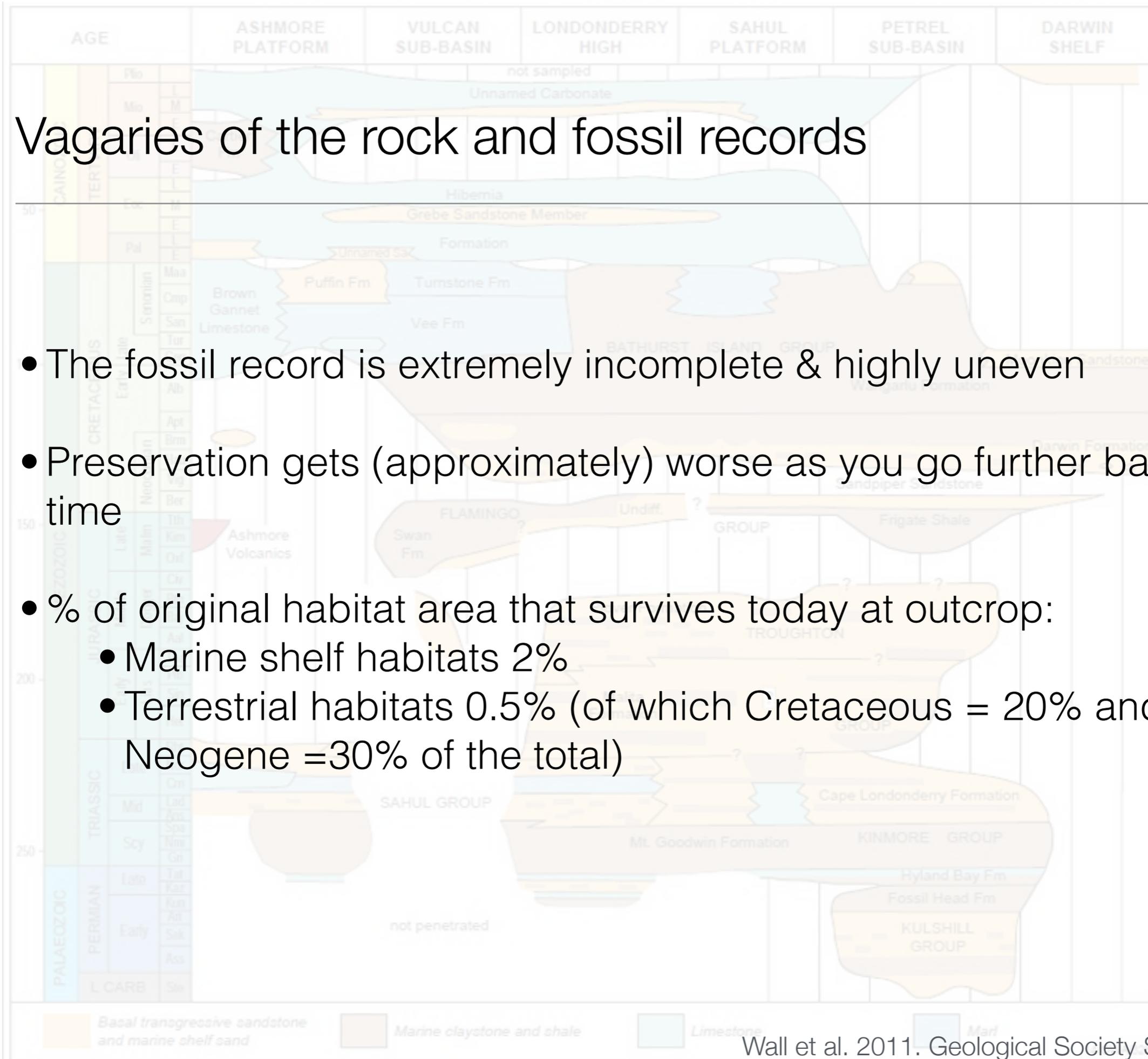
Vagaries of the rock and fossil records



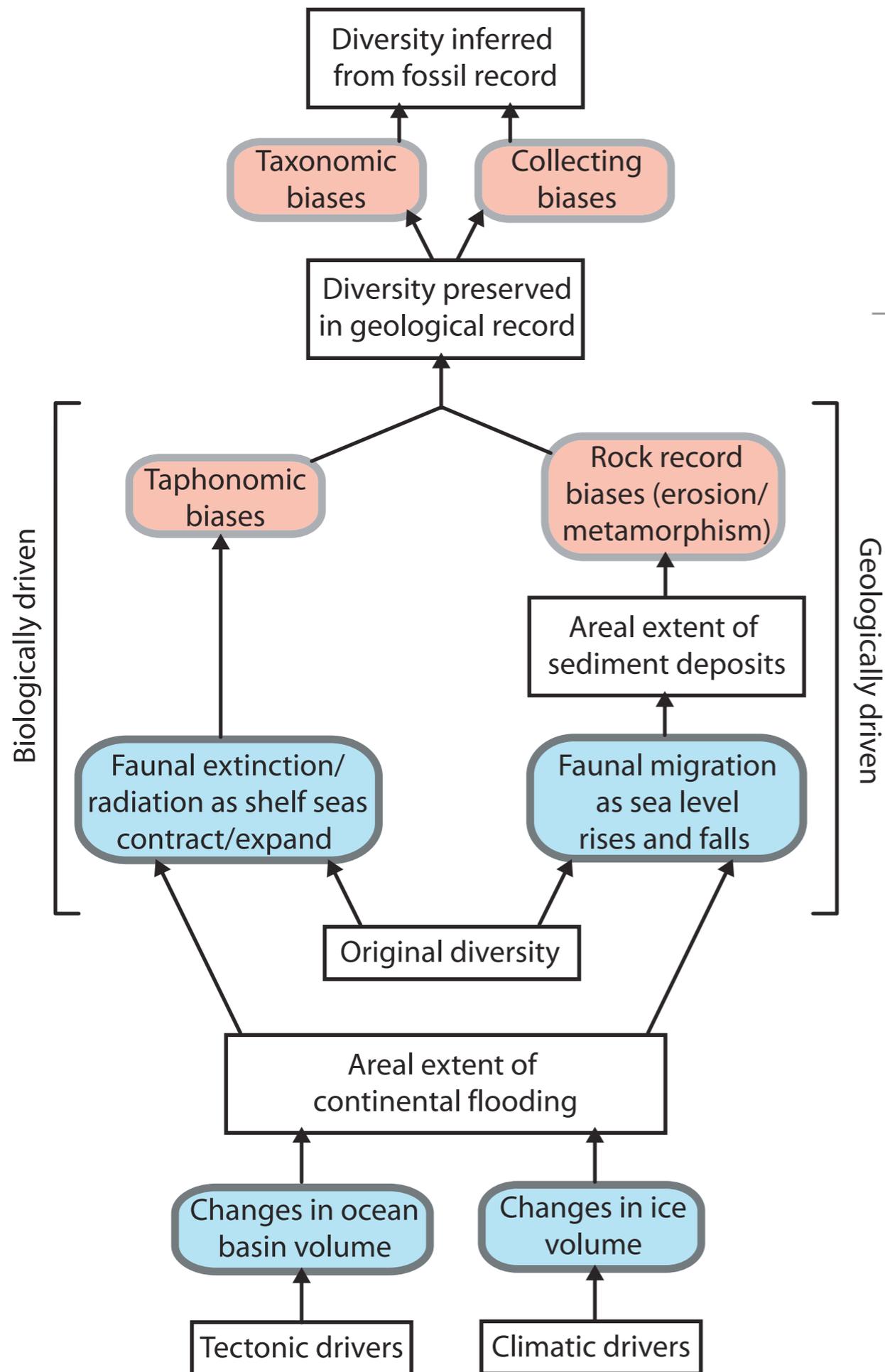


Vagaries of the rock and fossil records

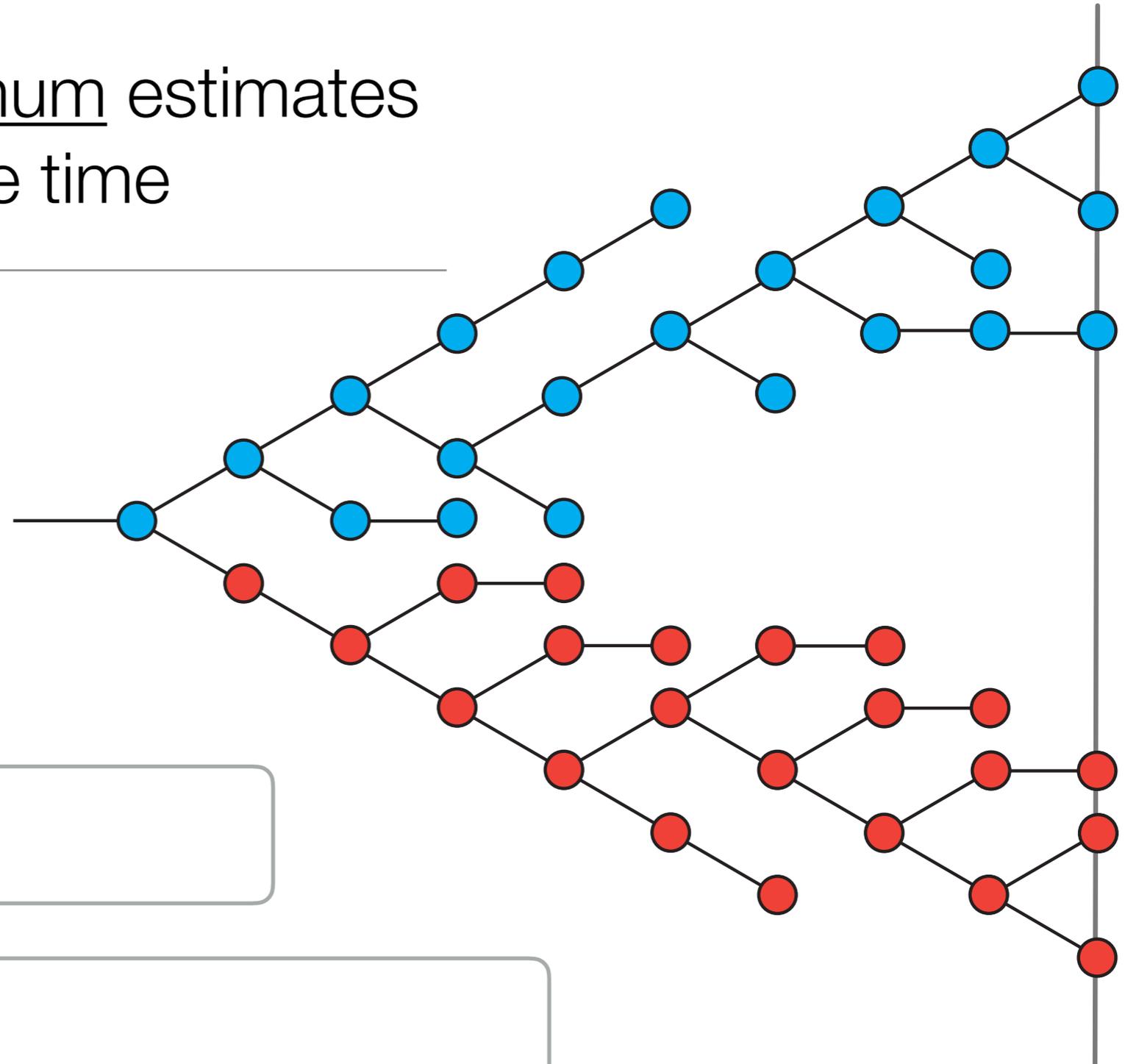
- The fossil record is extremely incomplete & highly uneven
- Preservation gets (approximately) worse as you go further back in time
- % of original habitat area that survives today at outcrop:
 - Marine shelf habitats 2%
 - Terrestrial habitats 0.5% (of which Cretaceous = 20% and Neogene = 30% of the total)



Controls on the probability of preservation



Fossils provide minimum estimates of species divergence time



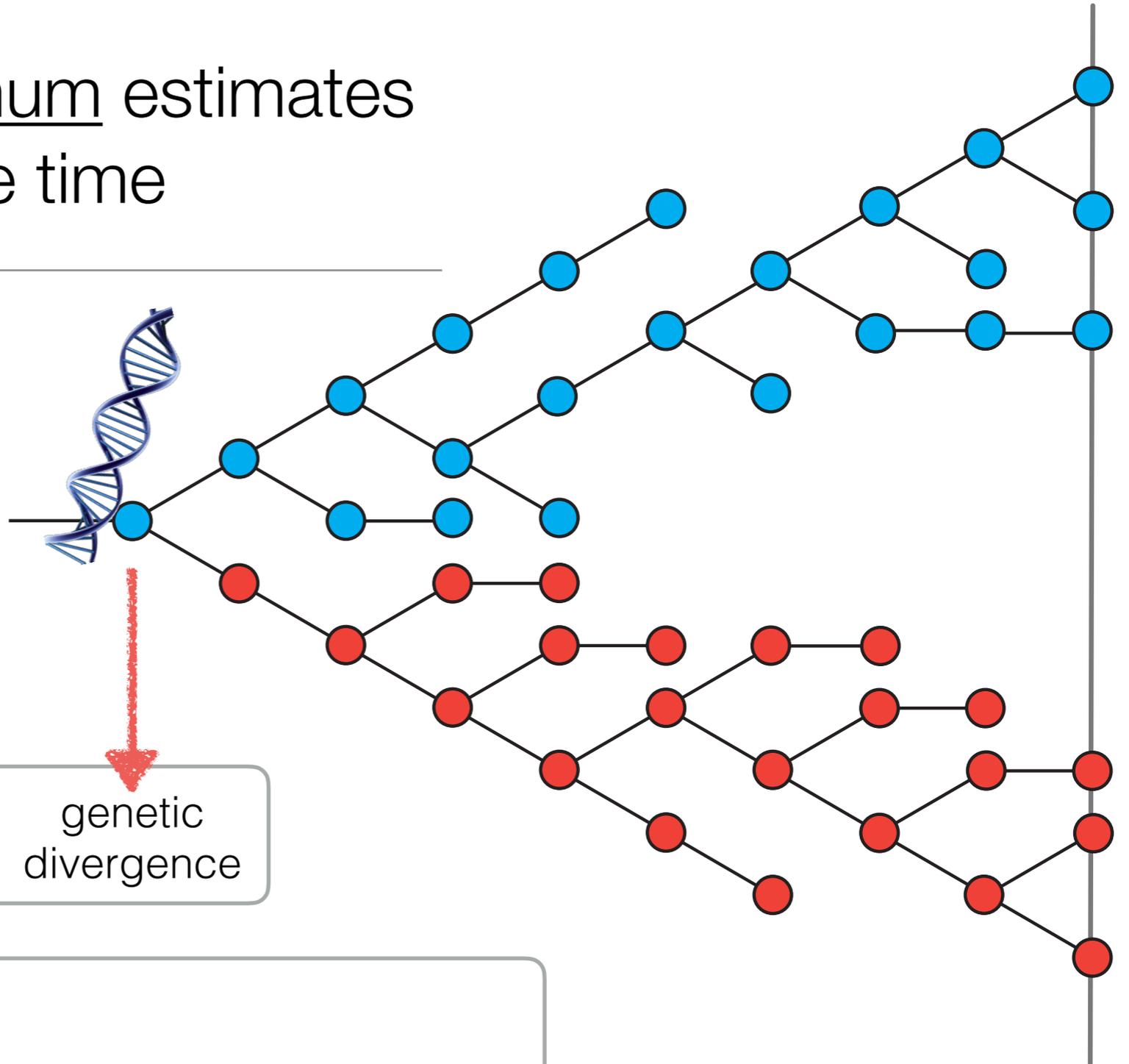
Molecular evolution:

Morphological evolution:

Fossil preservation:

Time

Fossils provide minimum estimates of species divergence time



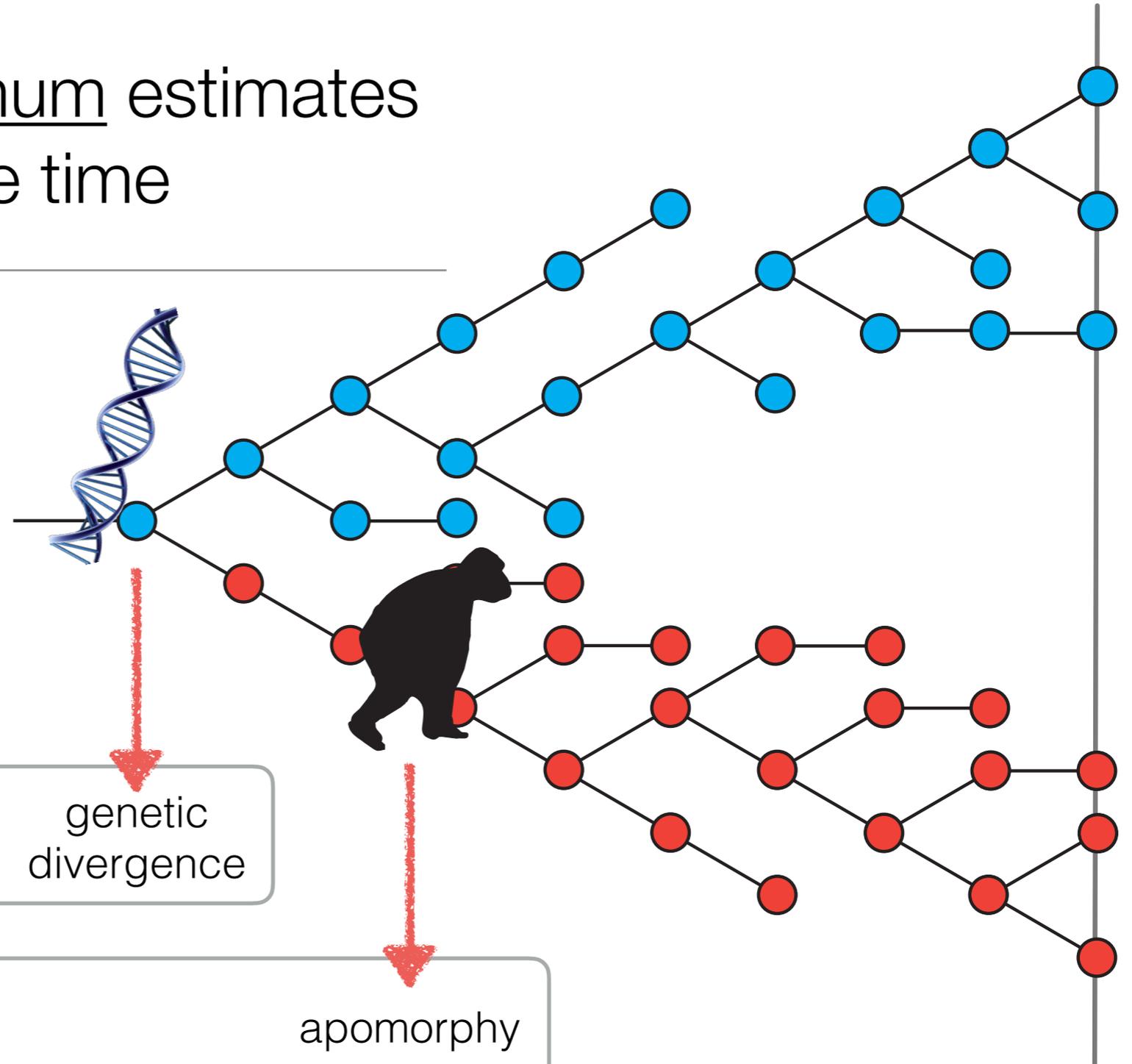
Molecular evolution: genetic divergence

Morphological evolution:

Fossil preservation:

Time

Fossils provide minimum estimates of species divergence time



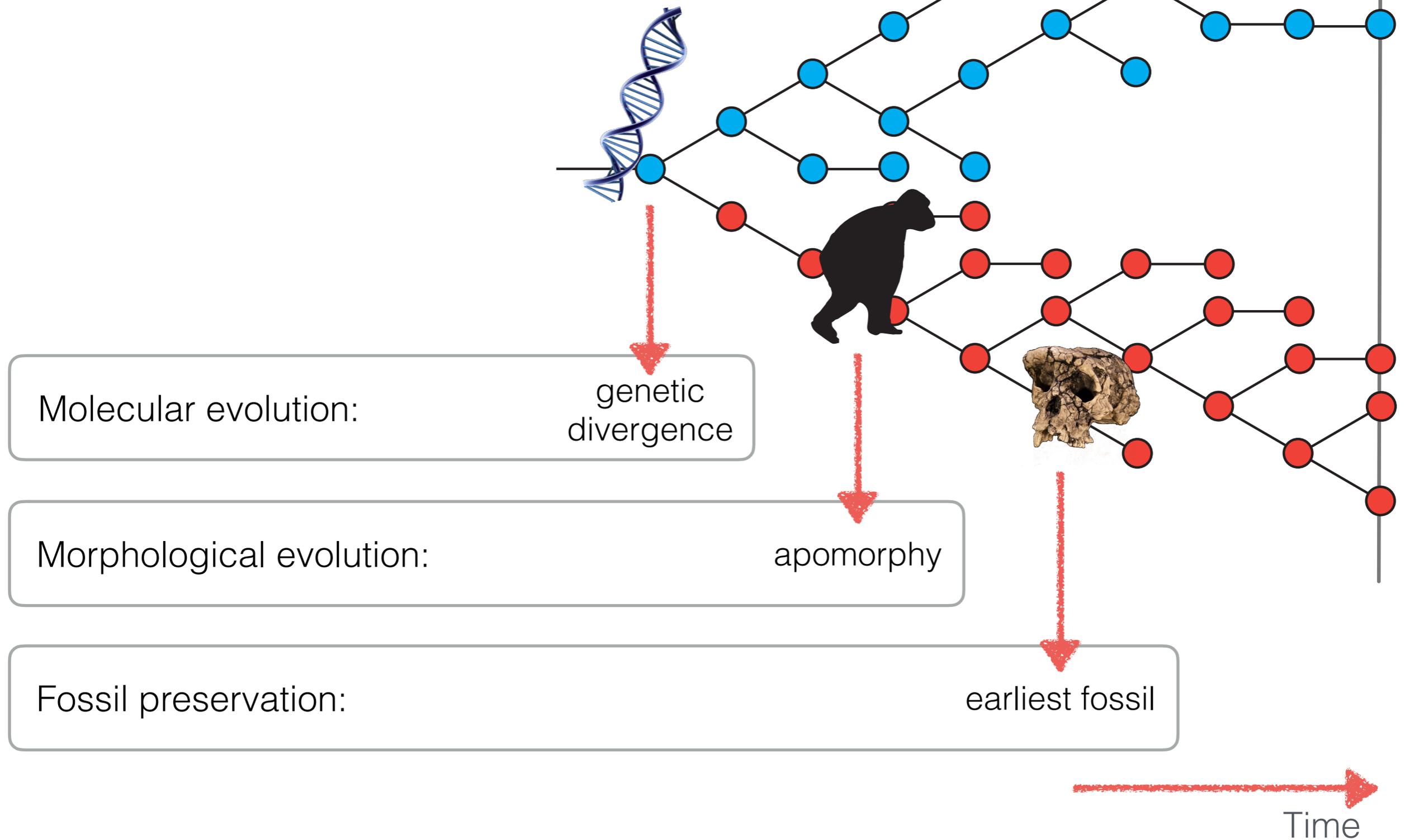
Molecular evolution: genetic divergence

Morphological evolution: apomorphy

Fossil preservation:

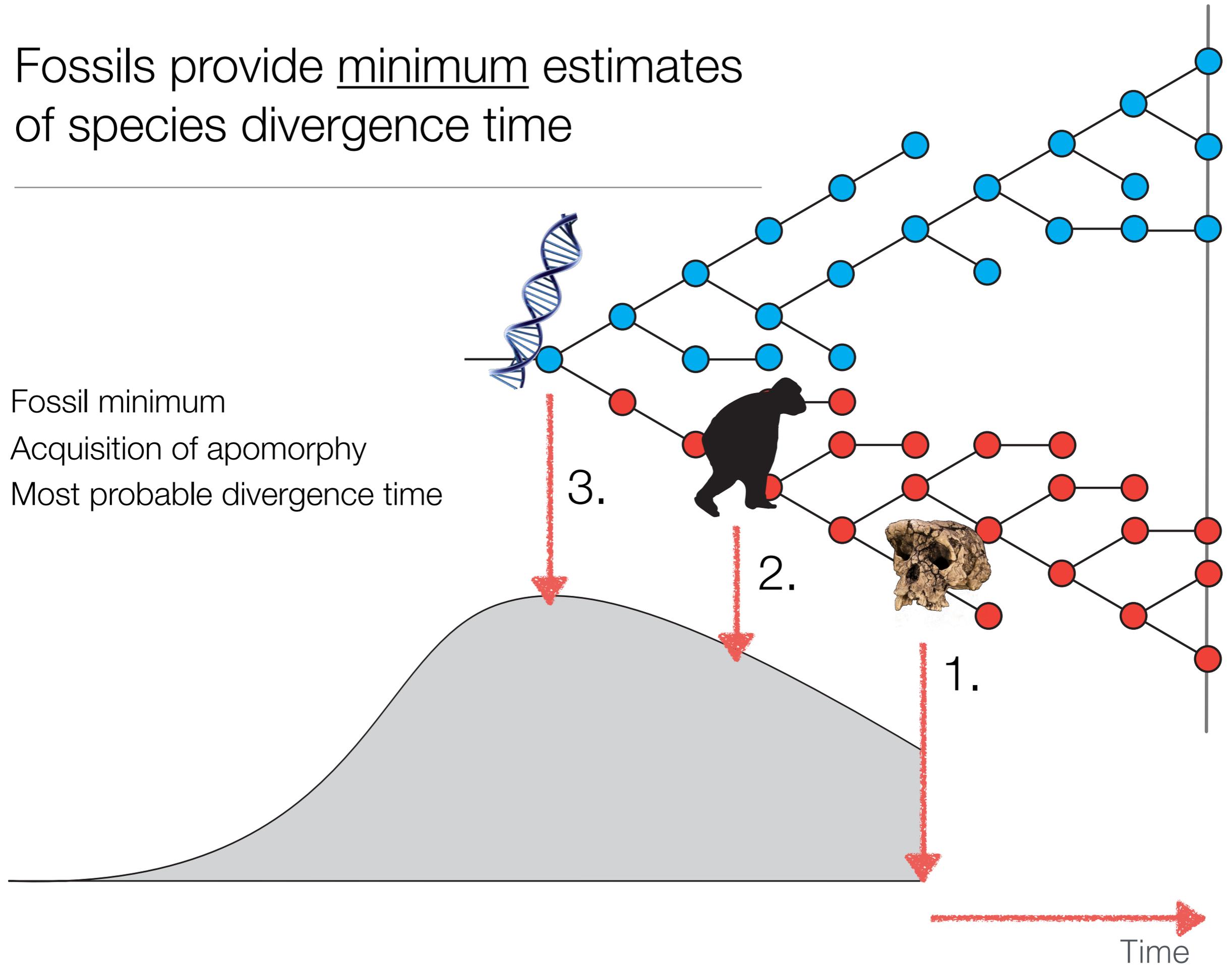
Time

Fossils provide minimum estimates of species divergence time

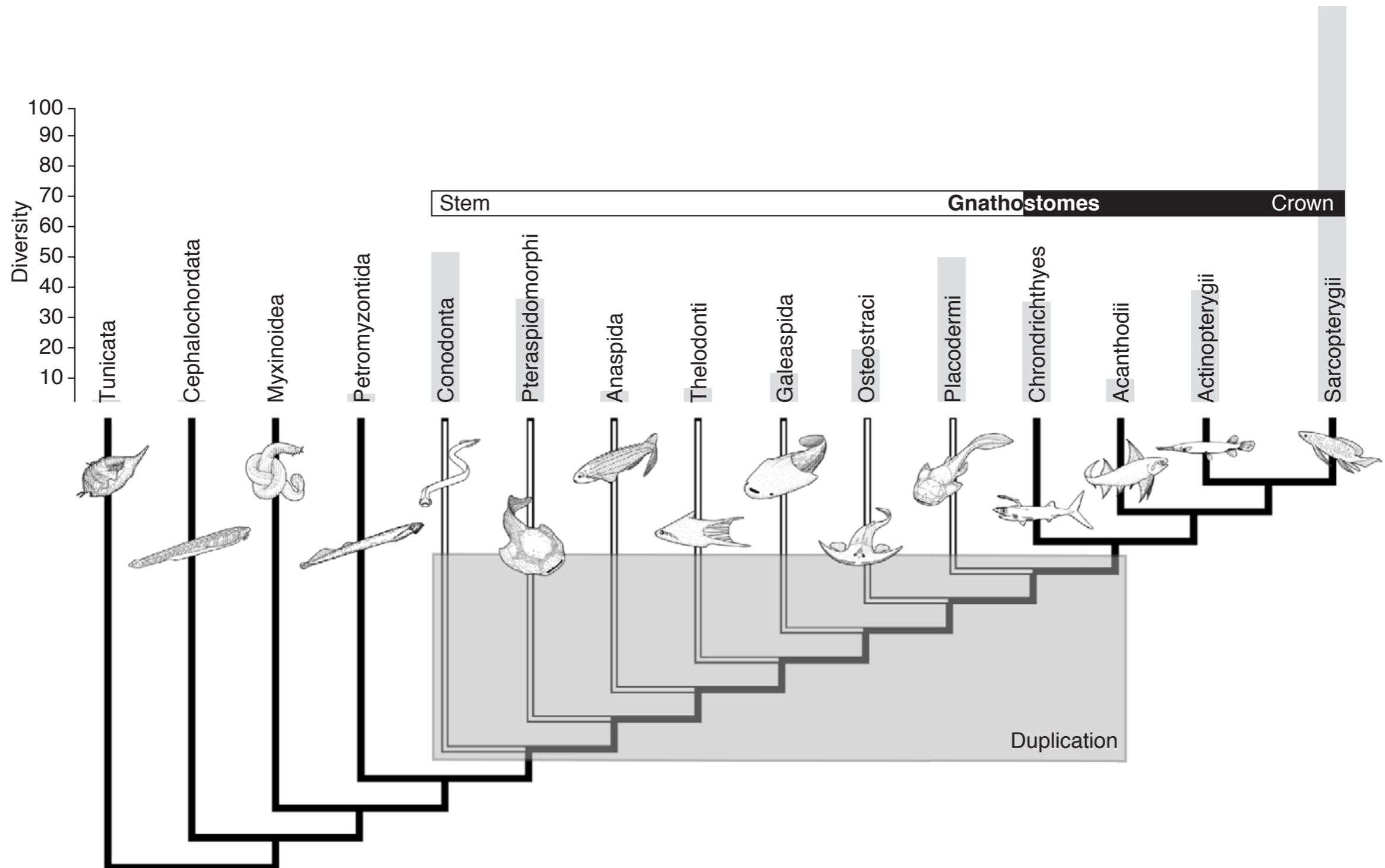


Fossils provide minimum estimates of species divergence time

1. Fossil minimum
2. Acquisition of apomorphy
3. Most probable divergence time

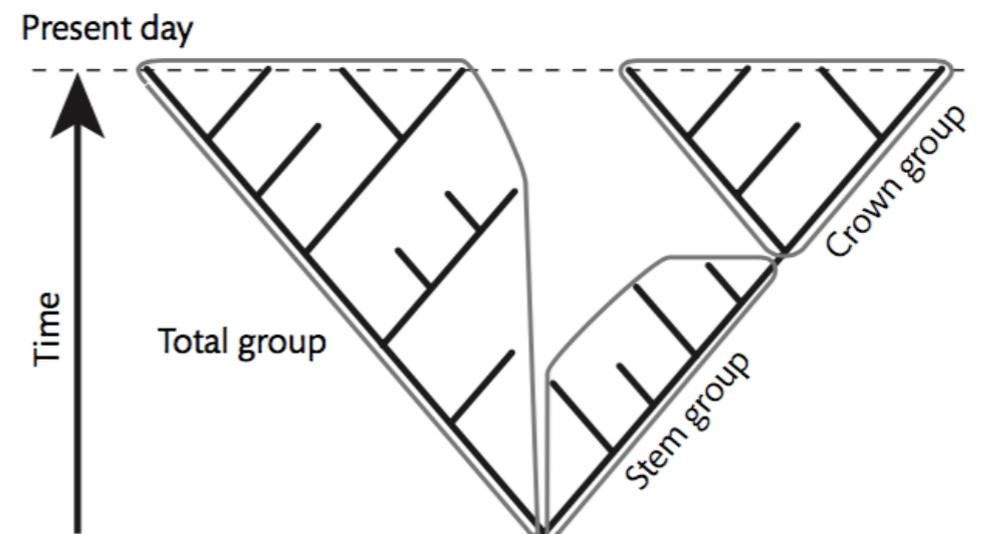


Taxonomic uncertainty: crown versus stem groups

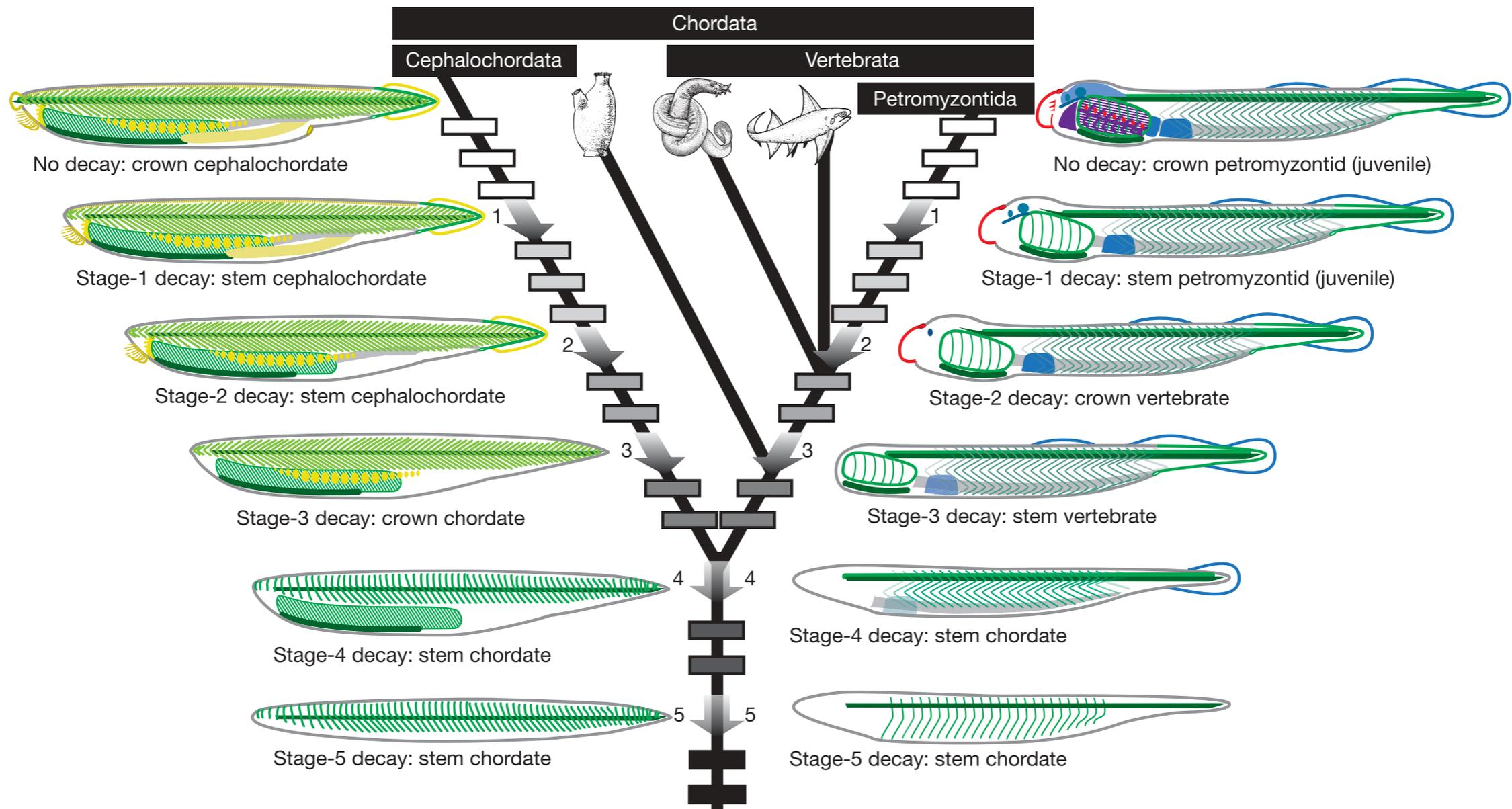


Taxonomic uncertainty: crown versus stem groups

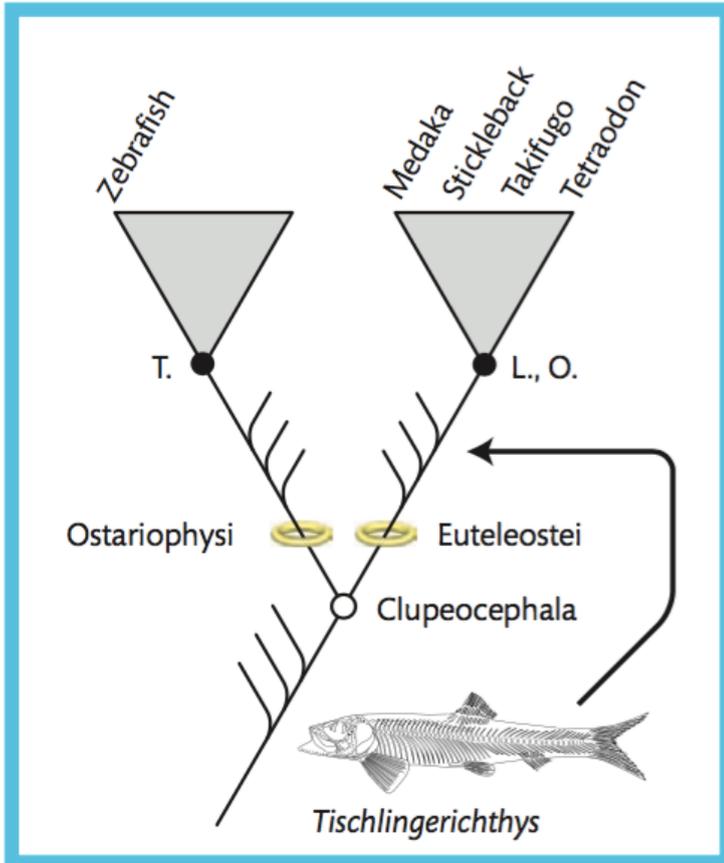
- Crown group: all descendants of the last common ancestor of the living members of a group
- Stem group: all species more closely related to the living members of a group than to any other
- Total group: stem + crown group members
- Early crown and stem group members may be difficult to distinguish



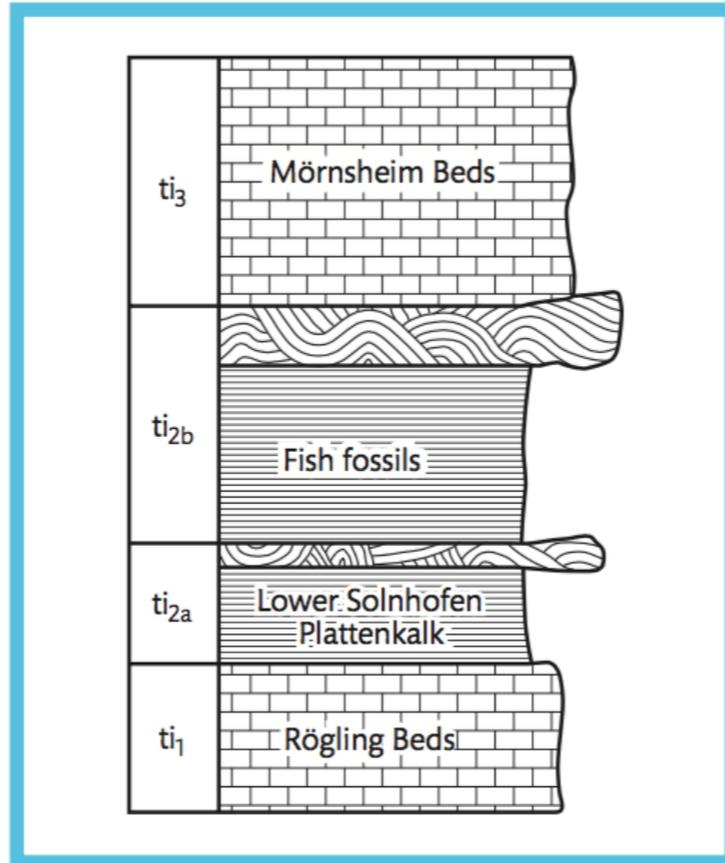
Taxonomic uncertainty: preservation biases



1. Oldest certain fossil in lineage

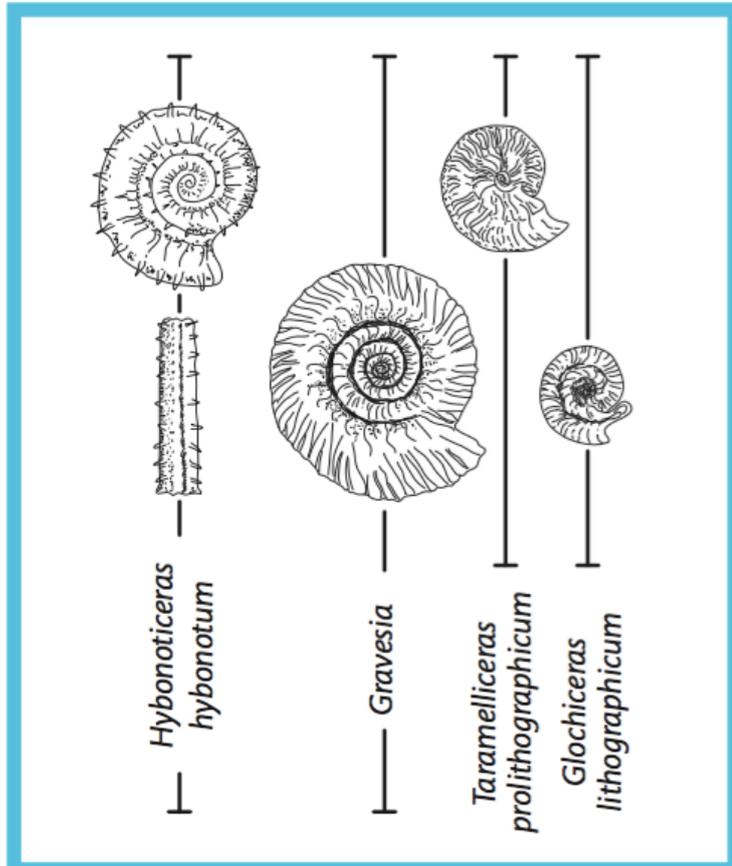


2. Lithostratigraphy of formation



Stratigraphic age uncertainty

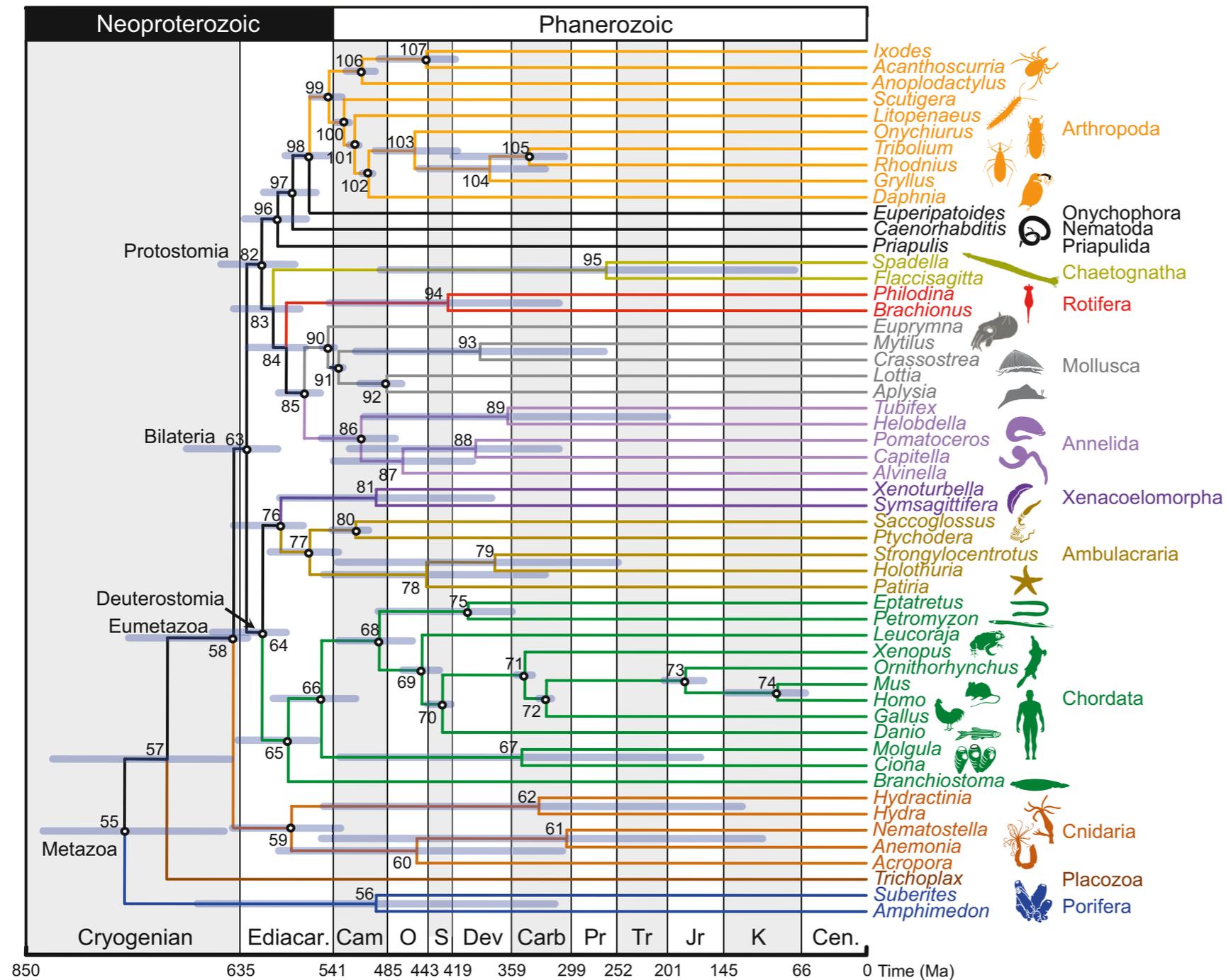
3. Biostratigraphy



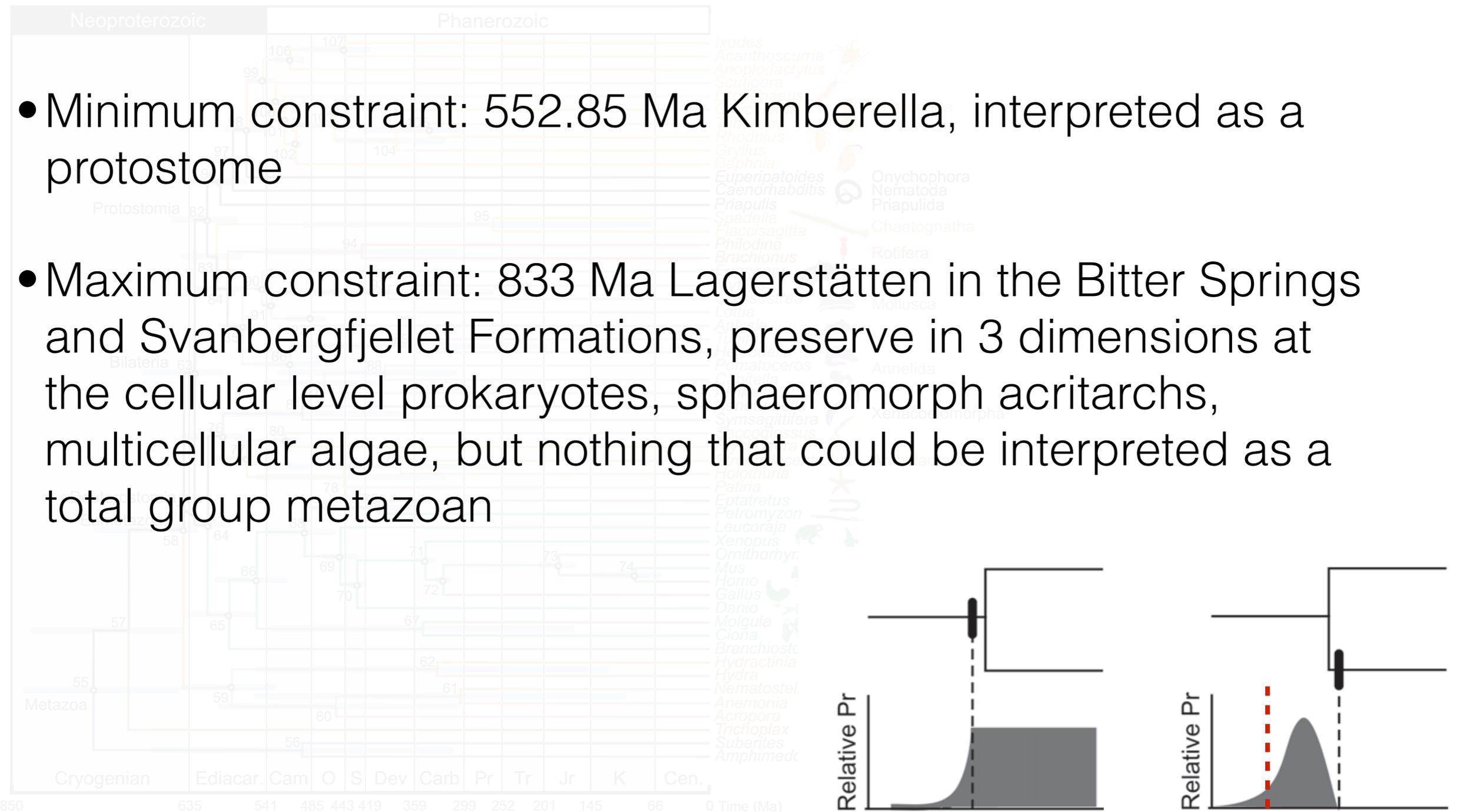
4. Chronostratigraphy

AGE (Ma)	Stage	Polarity Chron	Ammonite zones Sub-Mediterranean
145	CRETACEOUS (Berriasian)	M18	
145.5 ± 4.0		M19	
		M20	
	Tithonian	M21	
149.9 ± 0.05		M22	<i>S. semiforme</i> <i>S. darwini</i> <i>Hybonotoceras hybonotum</i>
150.8 ± 4.0		M22A	
		M23	<i>Hybonotoceras beckeri</i>
	Kimmeridgian	M24	
		M24A	
		M24B	
		M25	
154.55 ± 4.0		M25A	

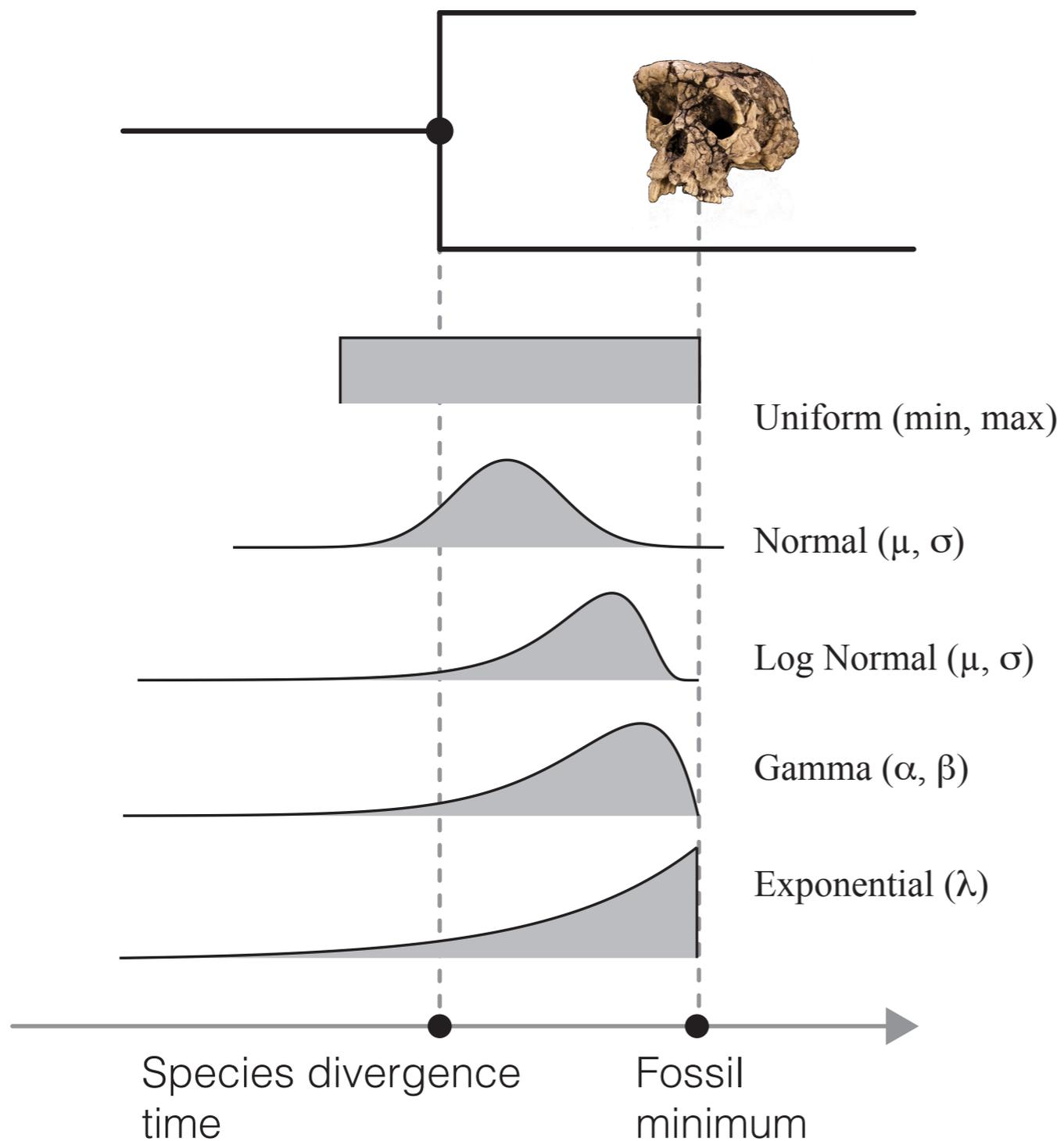
Soft maximum constraints on divergence times



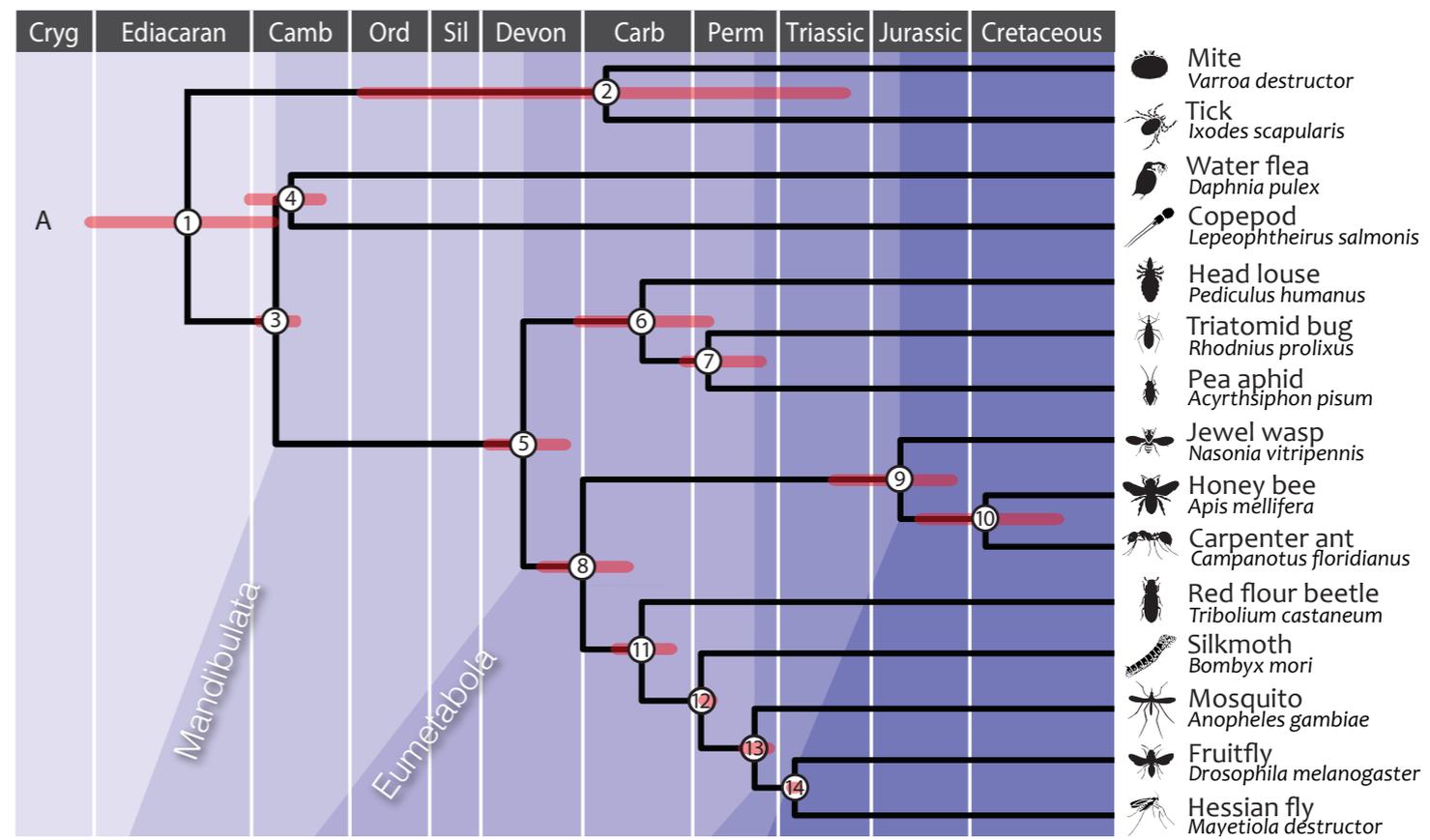
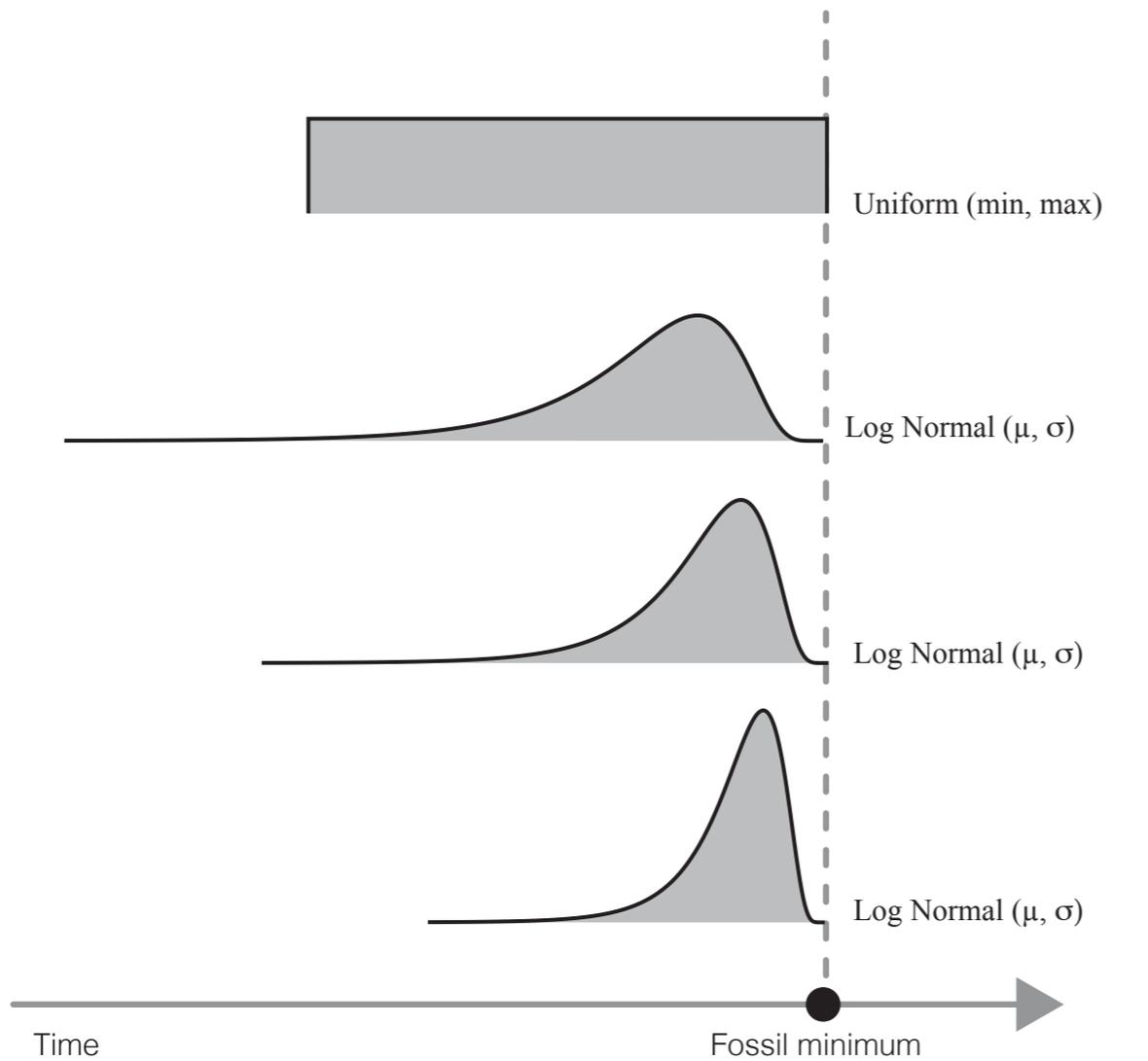
Soft maximum constraints on divergence times



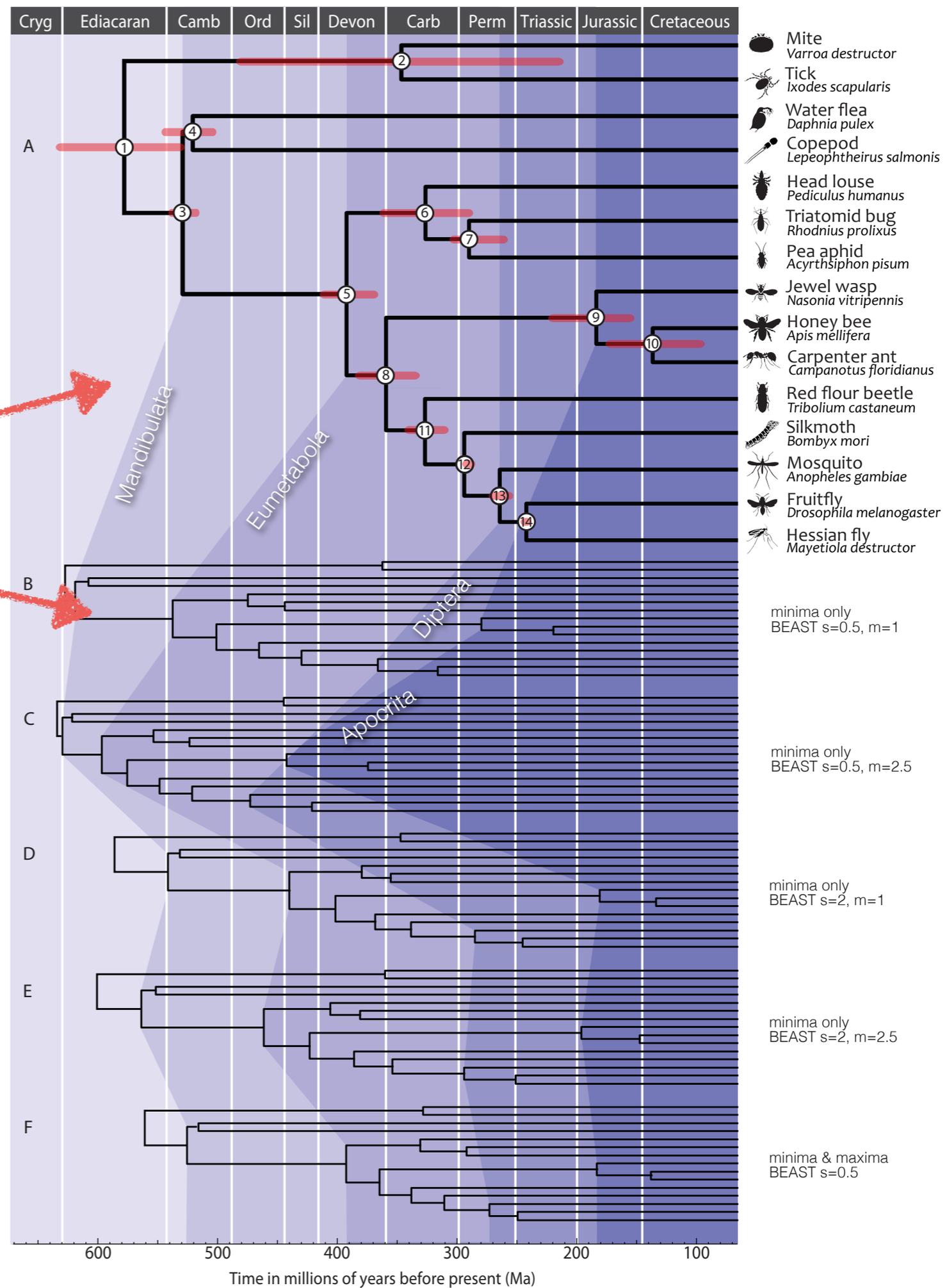
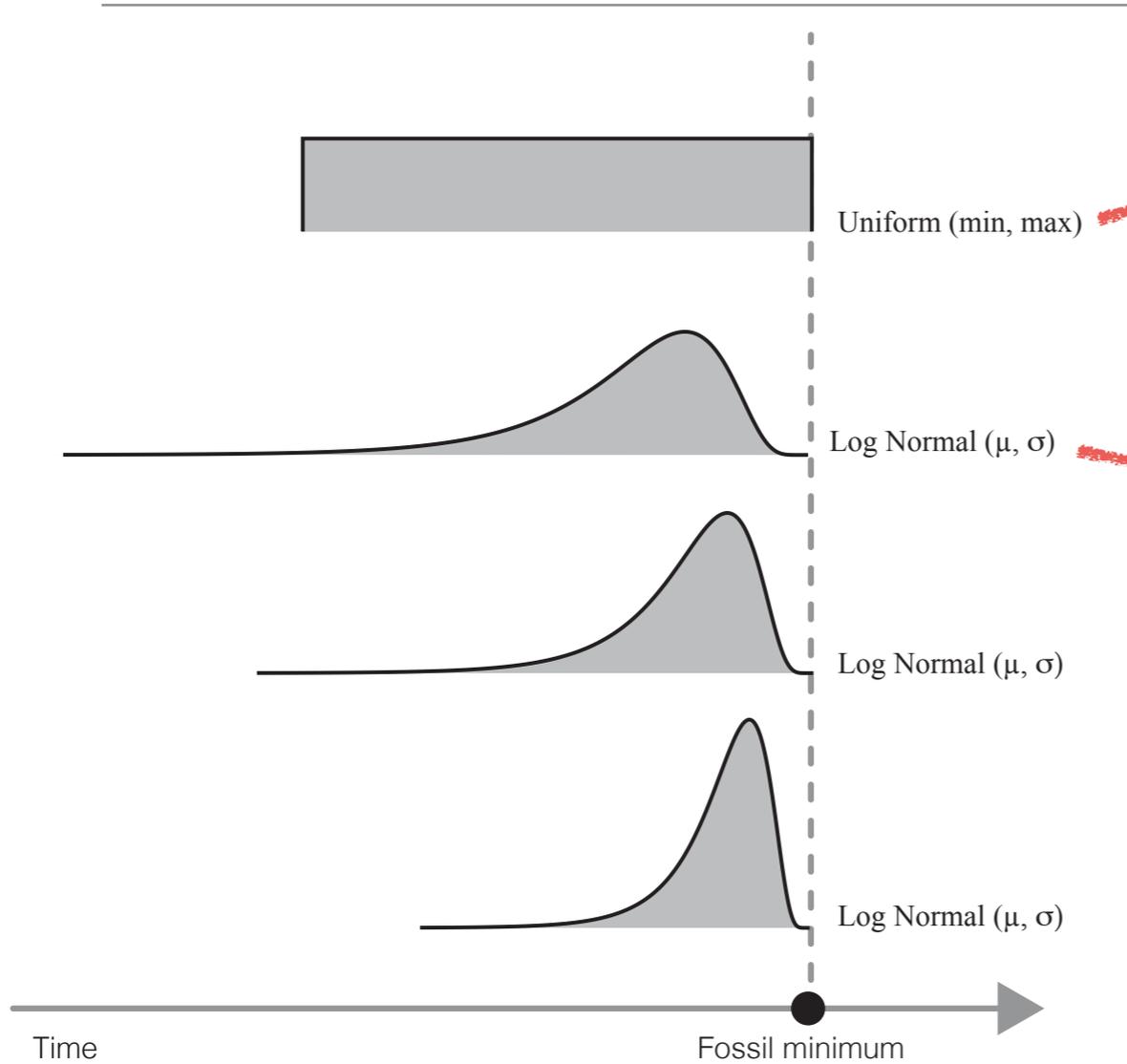
Probabilistic divergence times priors



The impact of different calibration priors



The impact of different calibration priors

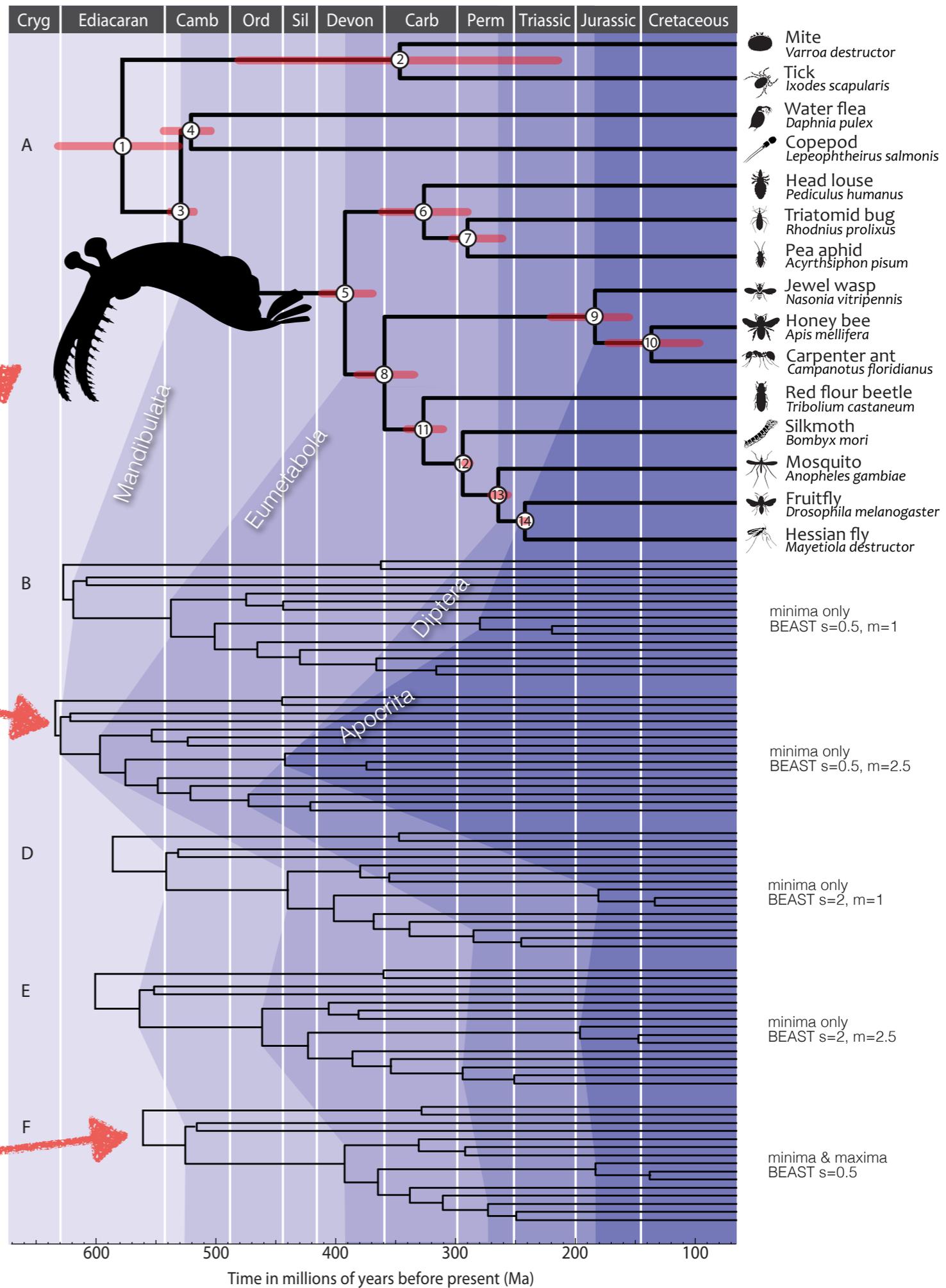


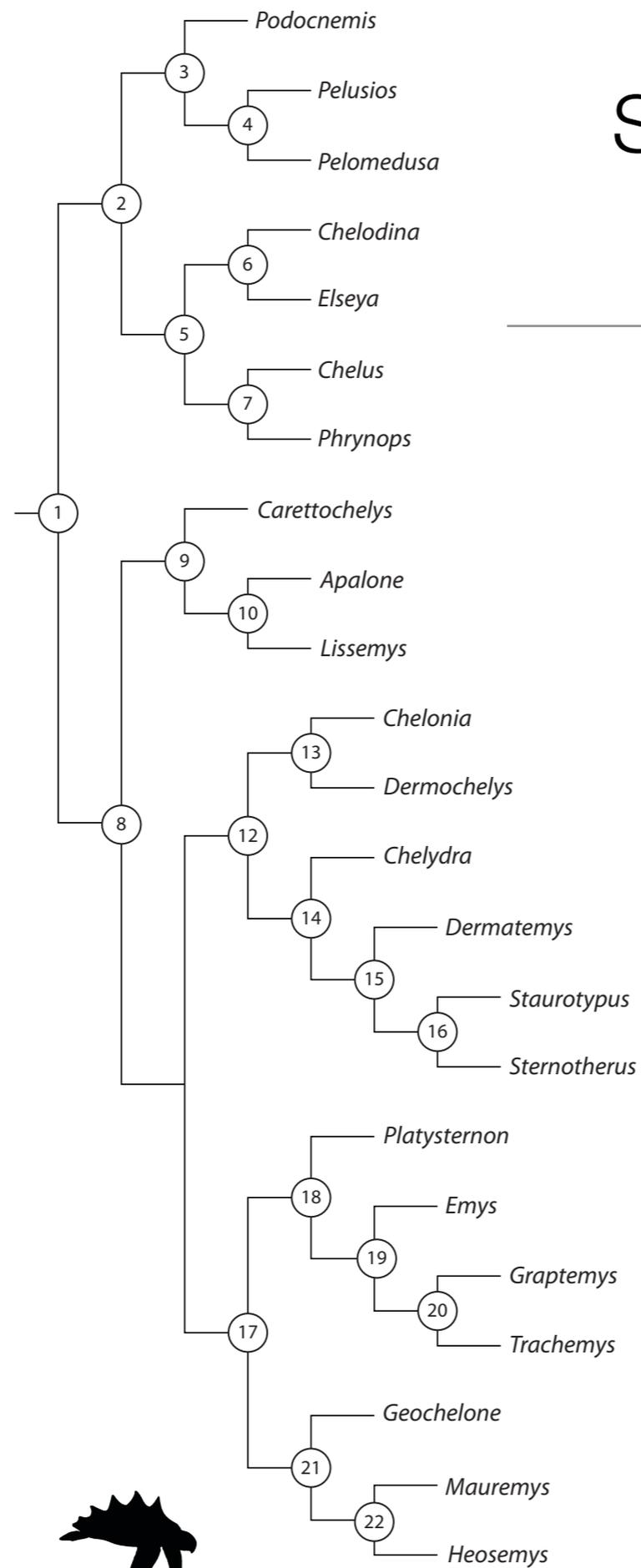
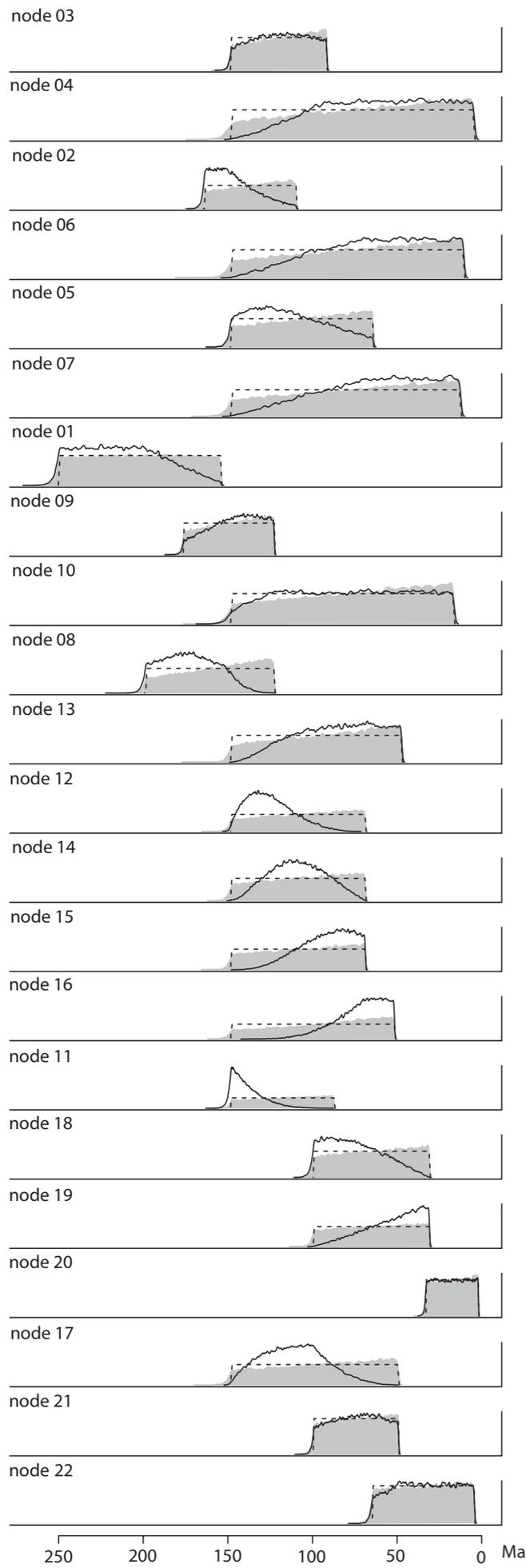
The impact of different calibration priors

Cambrian explosion

Cryogenian

Ediacaran



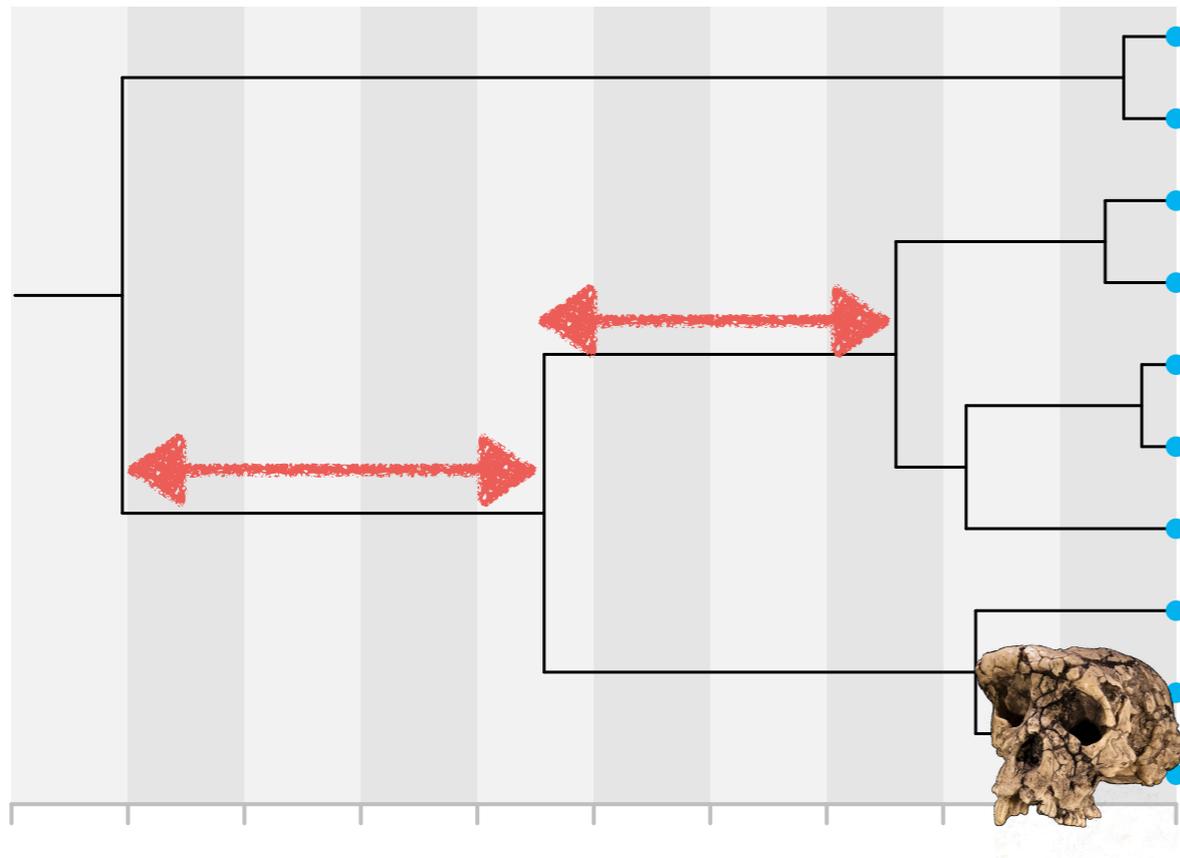


Specified versus effective priors



A prior for the non-fossil calibrated nodes

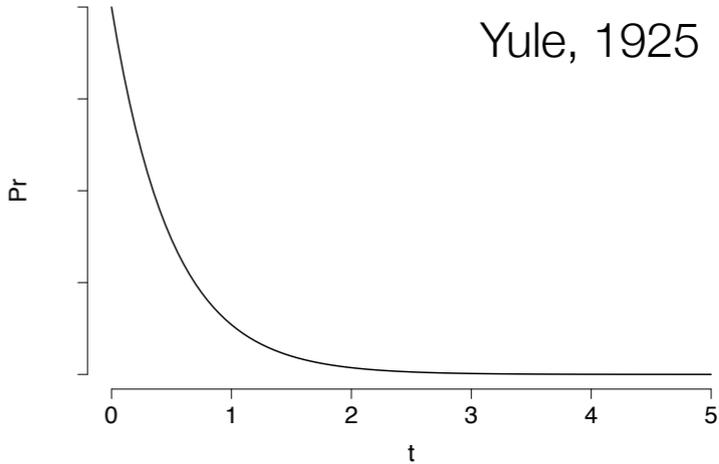
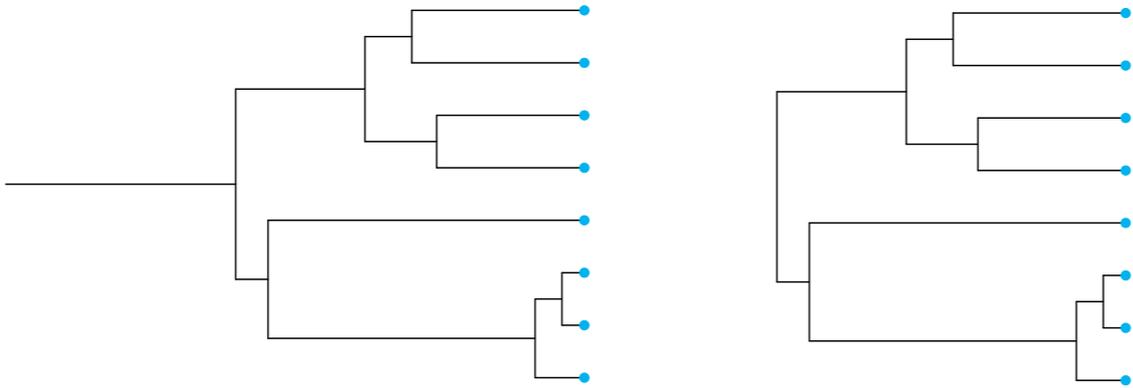
Extant species phylogeny



A prior for the uncalibrated nodes

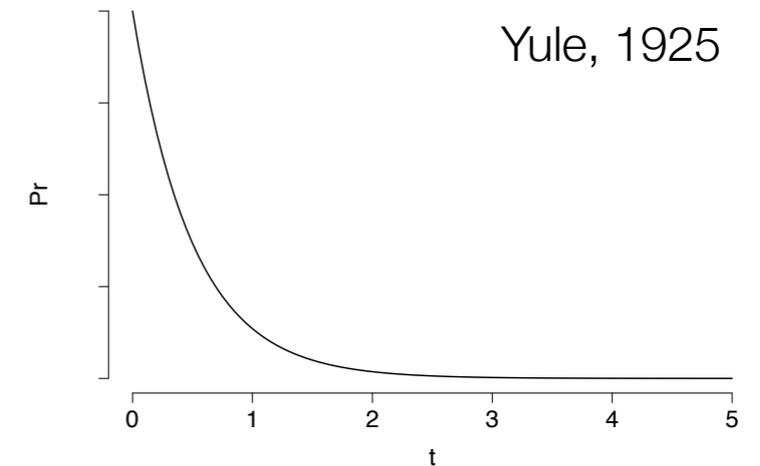
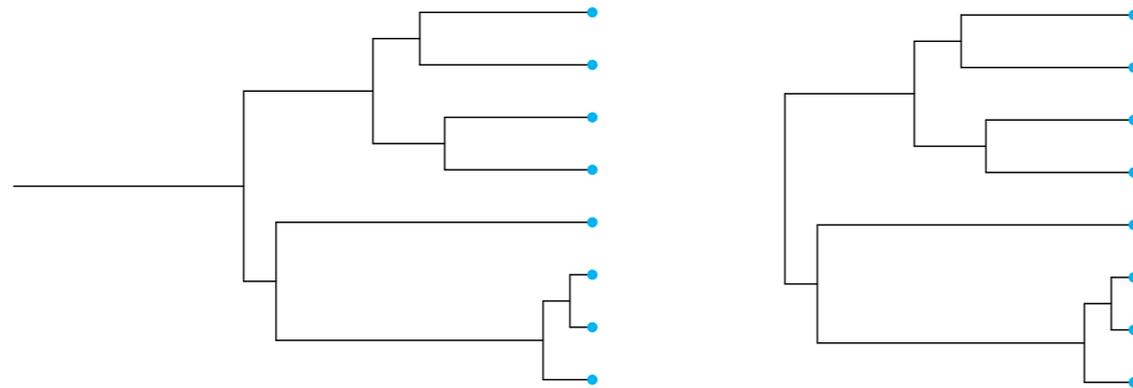
A prior for the non-fossil calibrated nodes: the tree prior

$\lambda = 1$

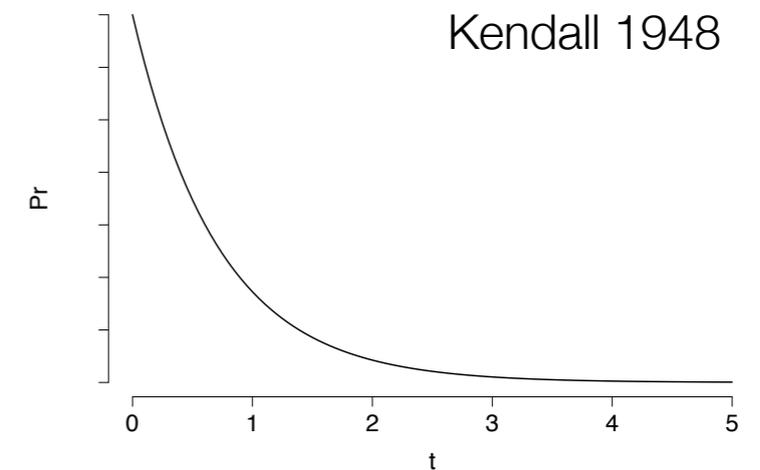
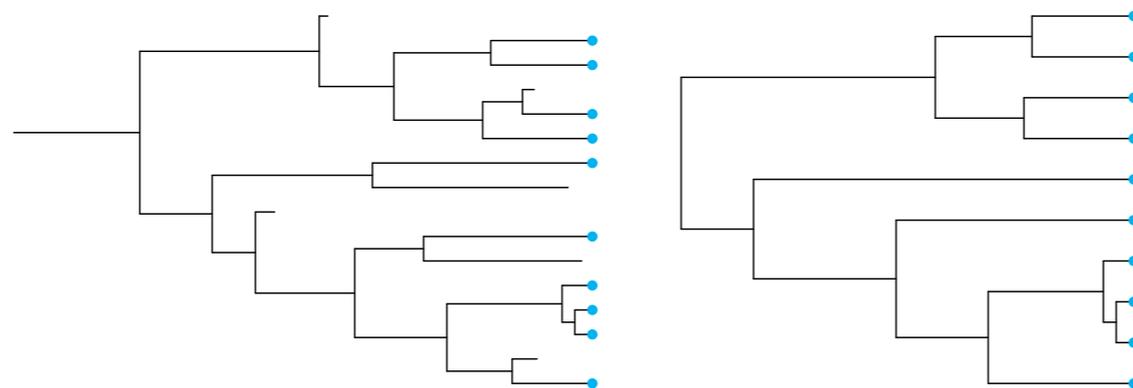


A prior for the non-fossil calibrated nodes: the tree prior

$\lambda = 1$



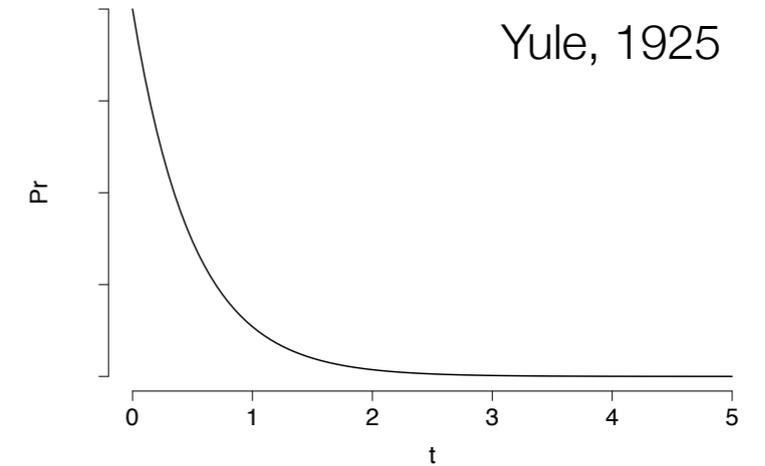
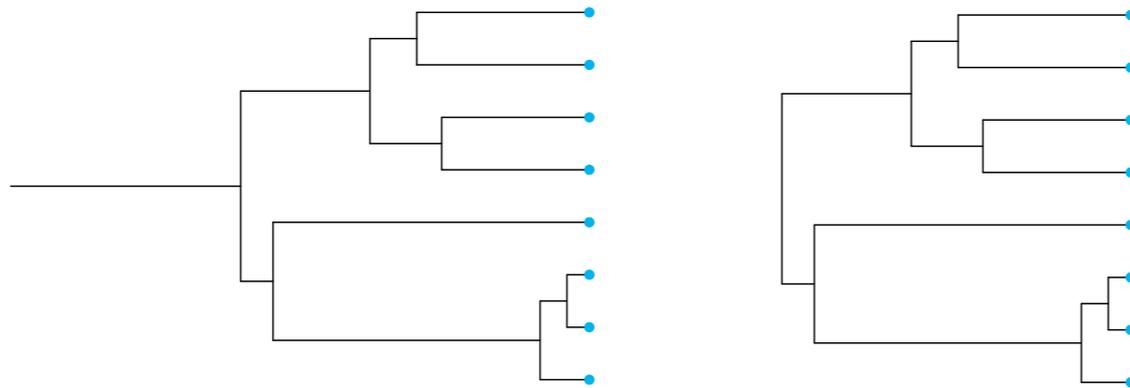
$\lambda = 1, \mu = 0.3$



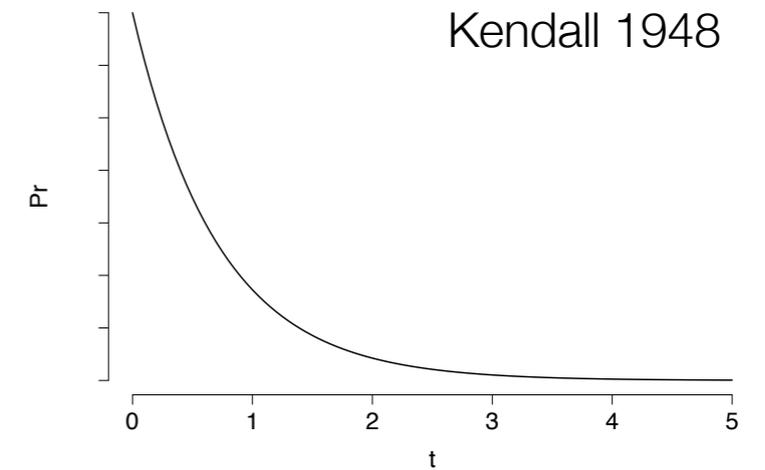
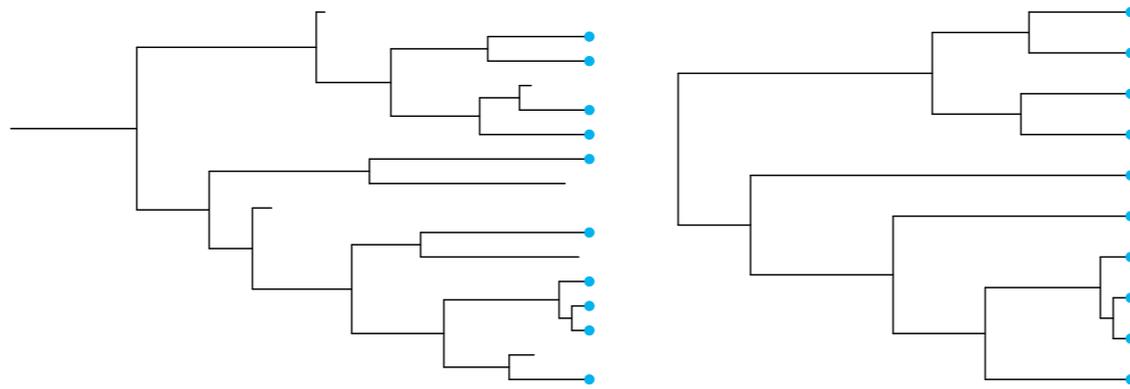
- Different combinations of λ and μ produce different tree shapes

A prior for the non-fossil calibrated nodes: the tree prior

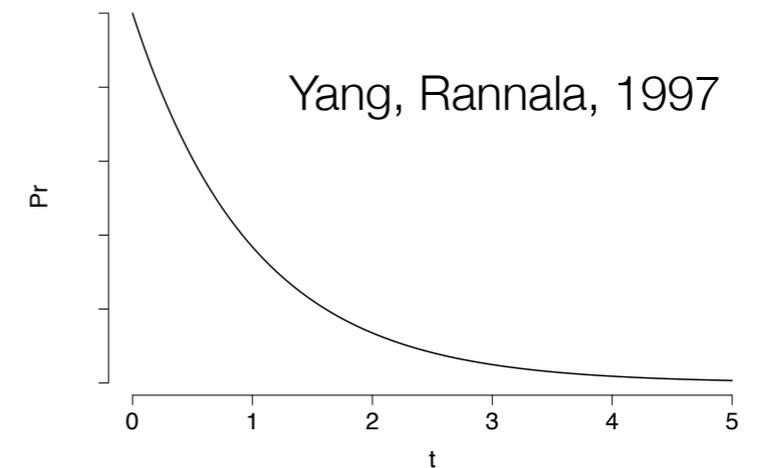
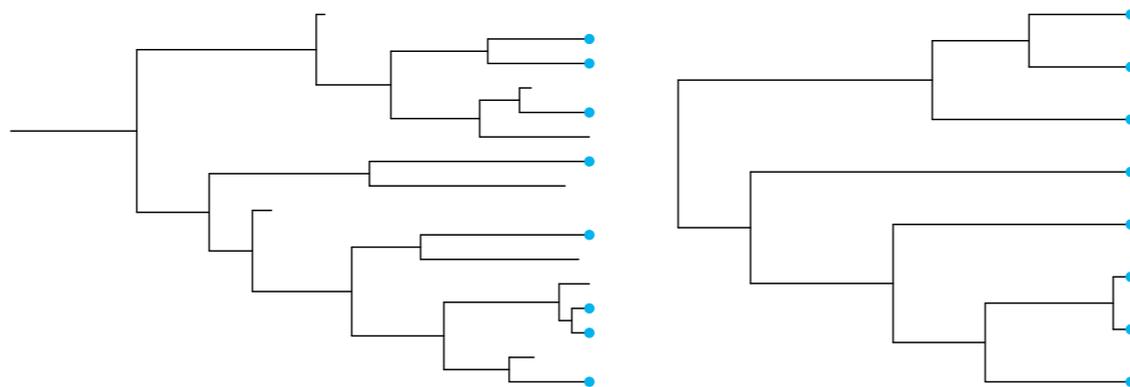
$\lambda = 1$



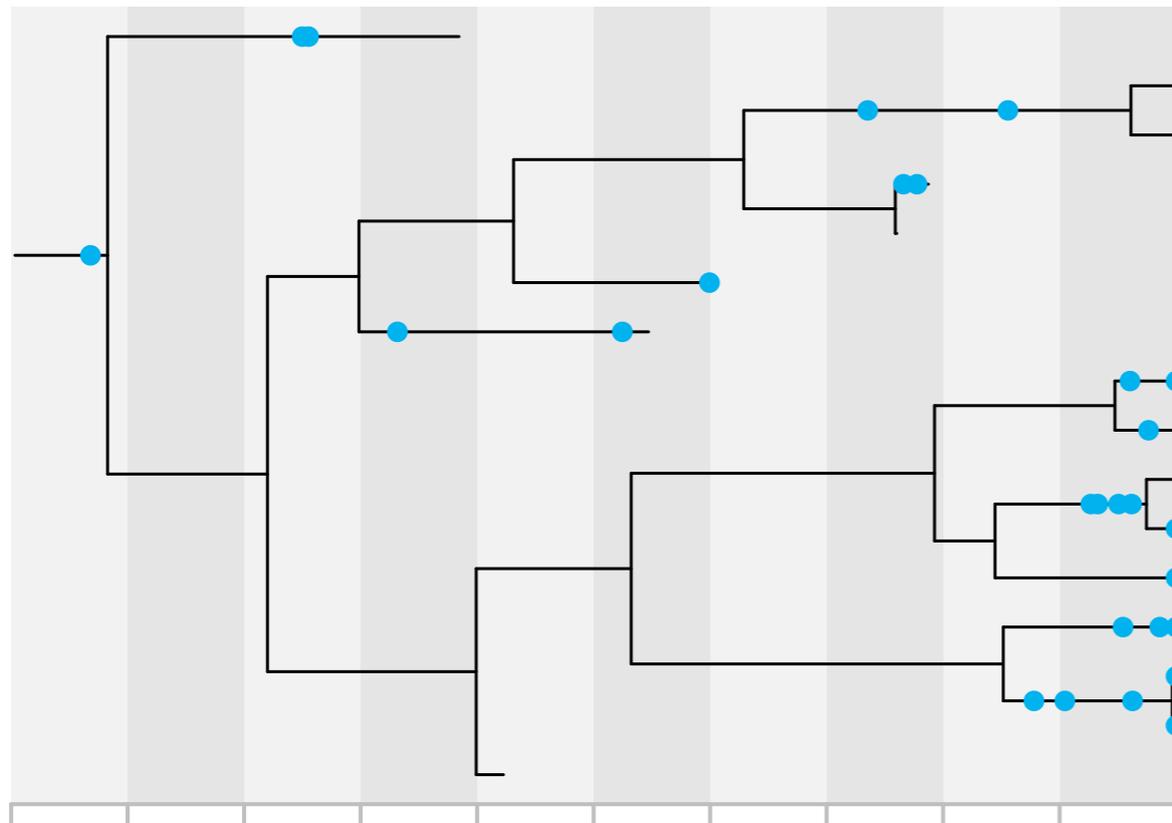
$\lambda = 1, \mu = 0.3$



$\lambda = 1, \mu = 0.3,$
 $\rho = 0.5$



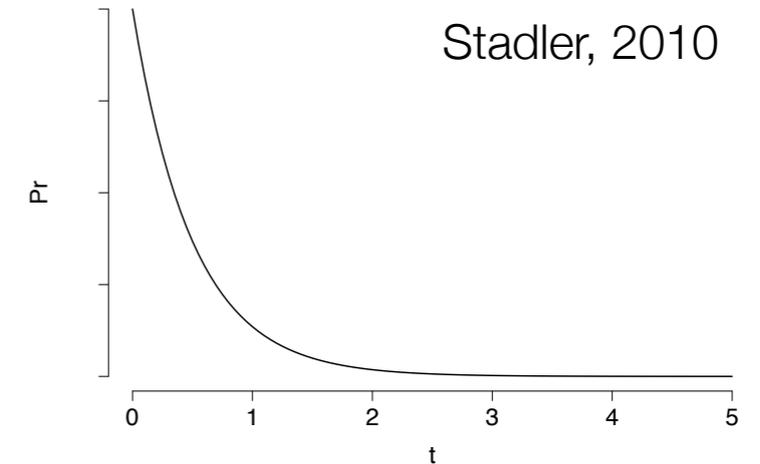
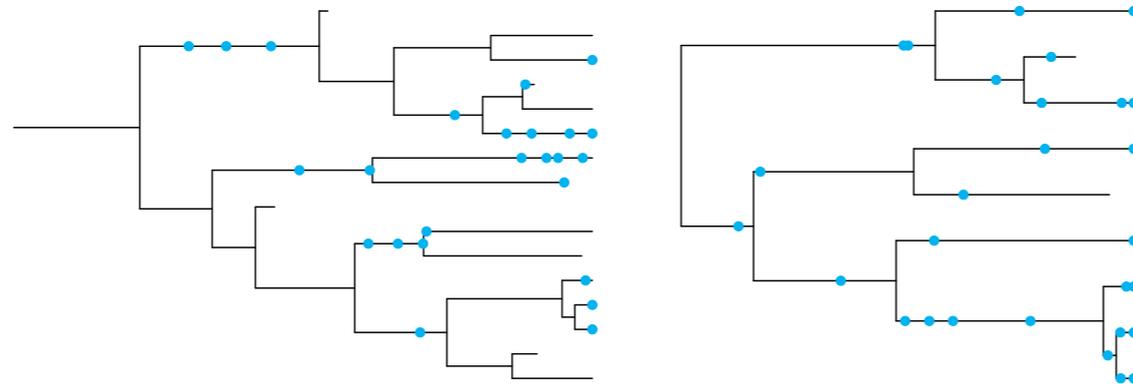
The fossilized birth-death process



- λ — speciation rate
- μ — extinction rate
- ψ — fossil sampling rate
- ρ — extant species sampling

Incorporating fossils into the tree prior

$$\lambda = 1, \mu = 0.3, \\ \rho = 0.5, \psi = 1$$



- The simpler models can be considered special cases of the FBD model
- Traditional node dating may depend on a small number of fossils
- Molecular clock analyses will be sensitive to both the calibration priors & the tree prior

Sampled ancestors

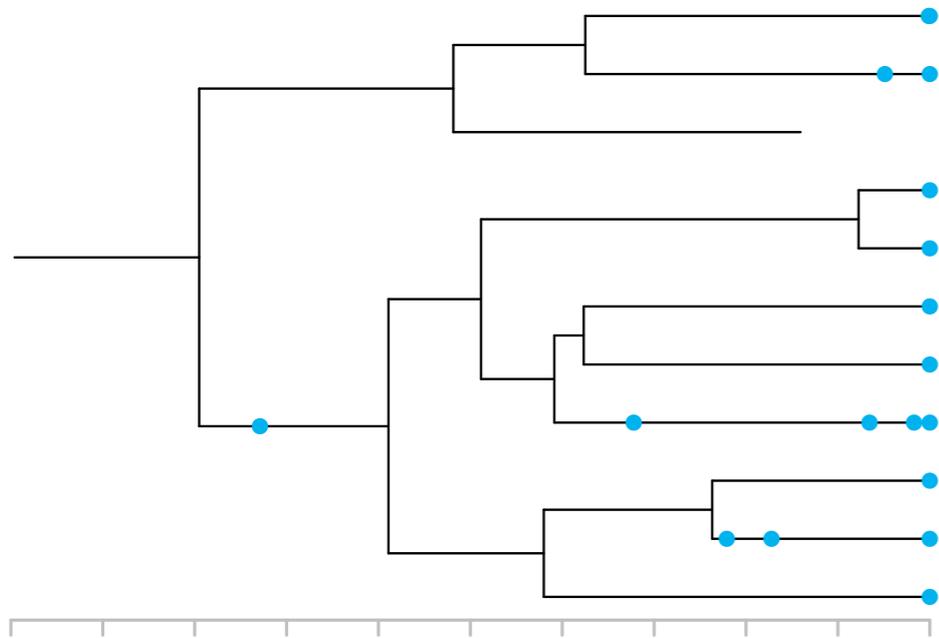
- The probability of sampling an ancestor in the fossil record is not zero (Foote, 1996. Paleobiology)

- Bayesian inference of sampled ancestor trees now possible in BEAST2, MrBayes, RevBayes, DPPDiv

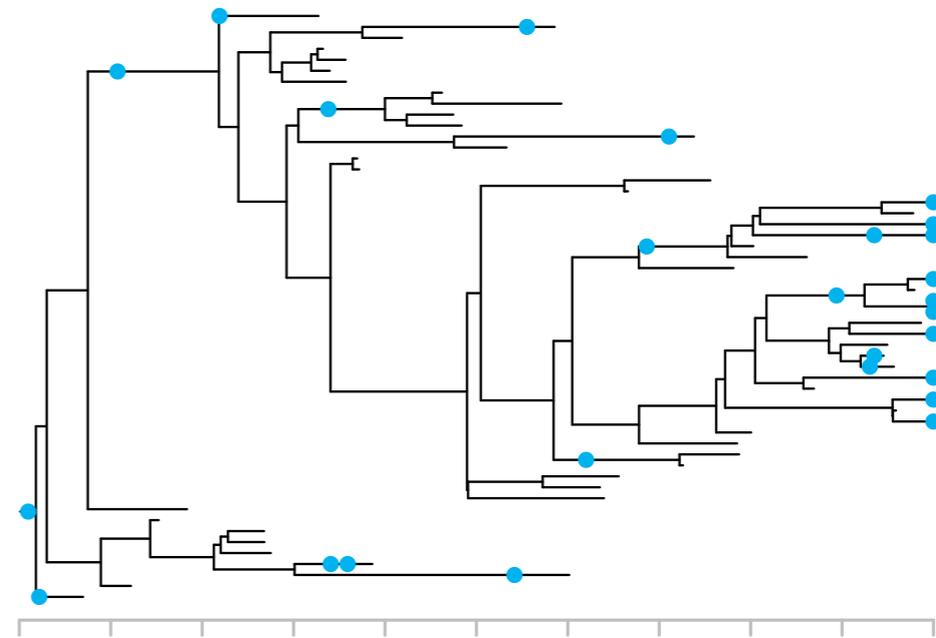
- Different combinations of λ , μ , ρ , ψ affect the probability of having a sampling an ancestor (Wright, Heath. in prep)

$\lambda = 0.1, \mu = 0.05, r = 0.5$

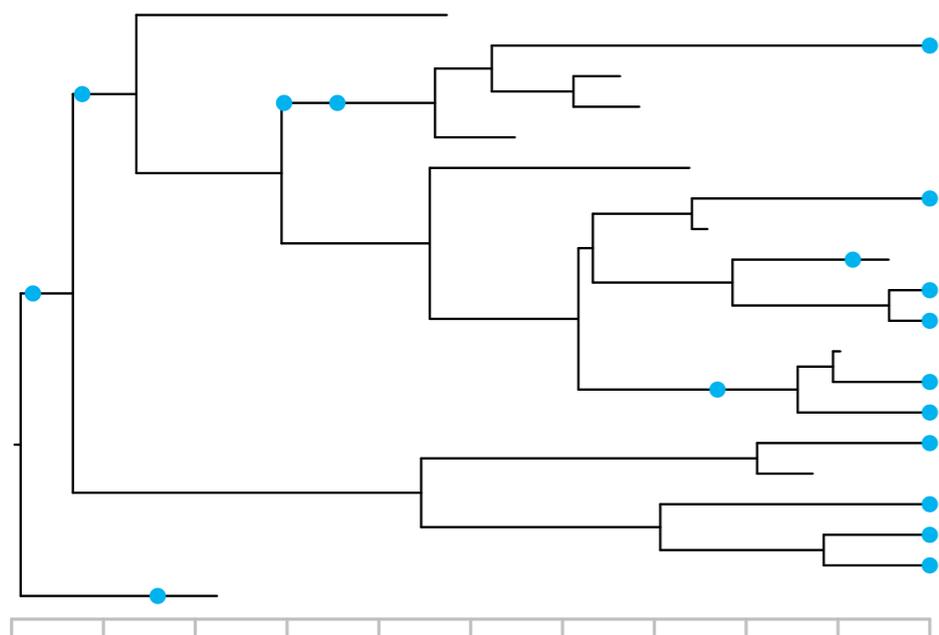
Sampled ancestors



$\lambda = 0.1, \mu = 0.03, r = 0.3$

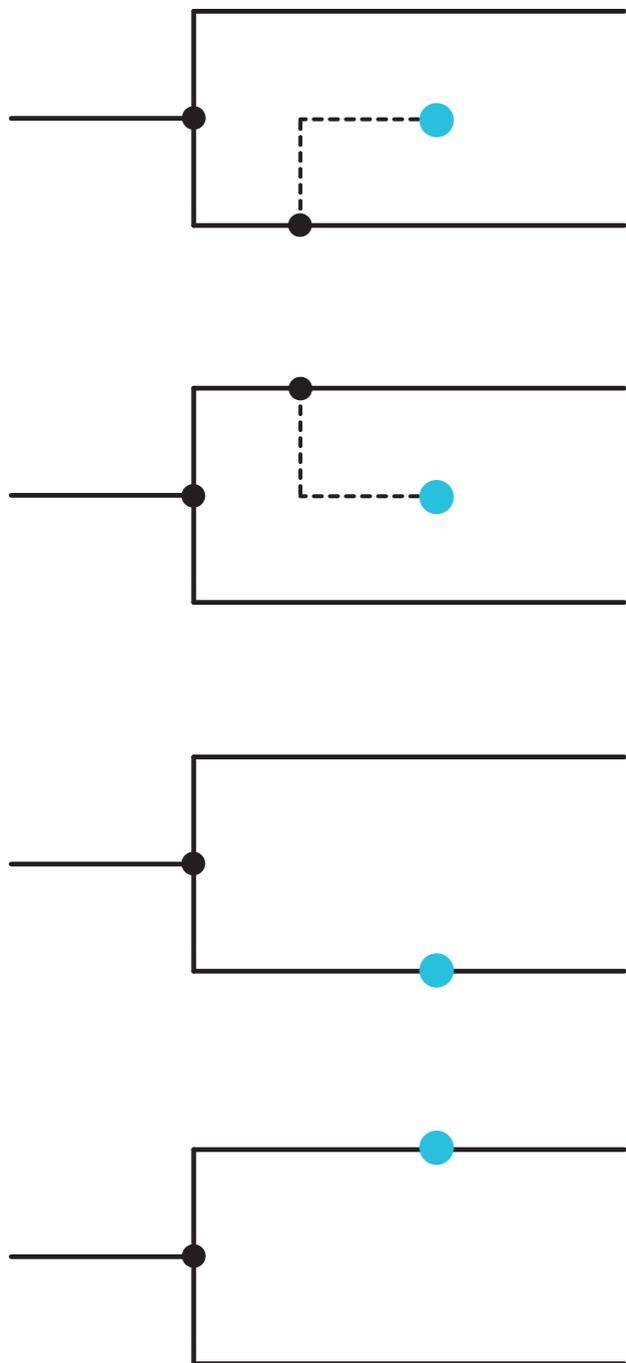


$\lambda = 0.1, \mu = 0.08, r = 0.8$



$\lambda = 0.1, \mu = 0.05, r = 0.5$

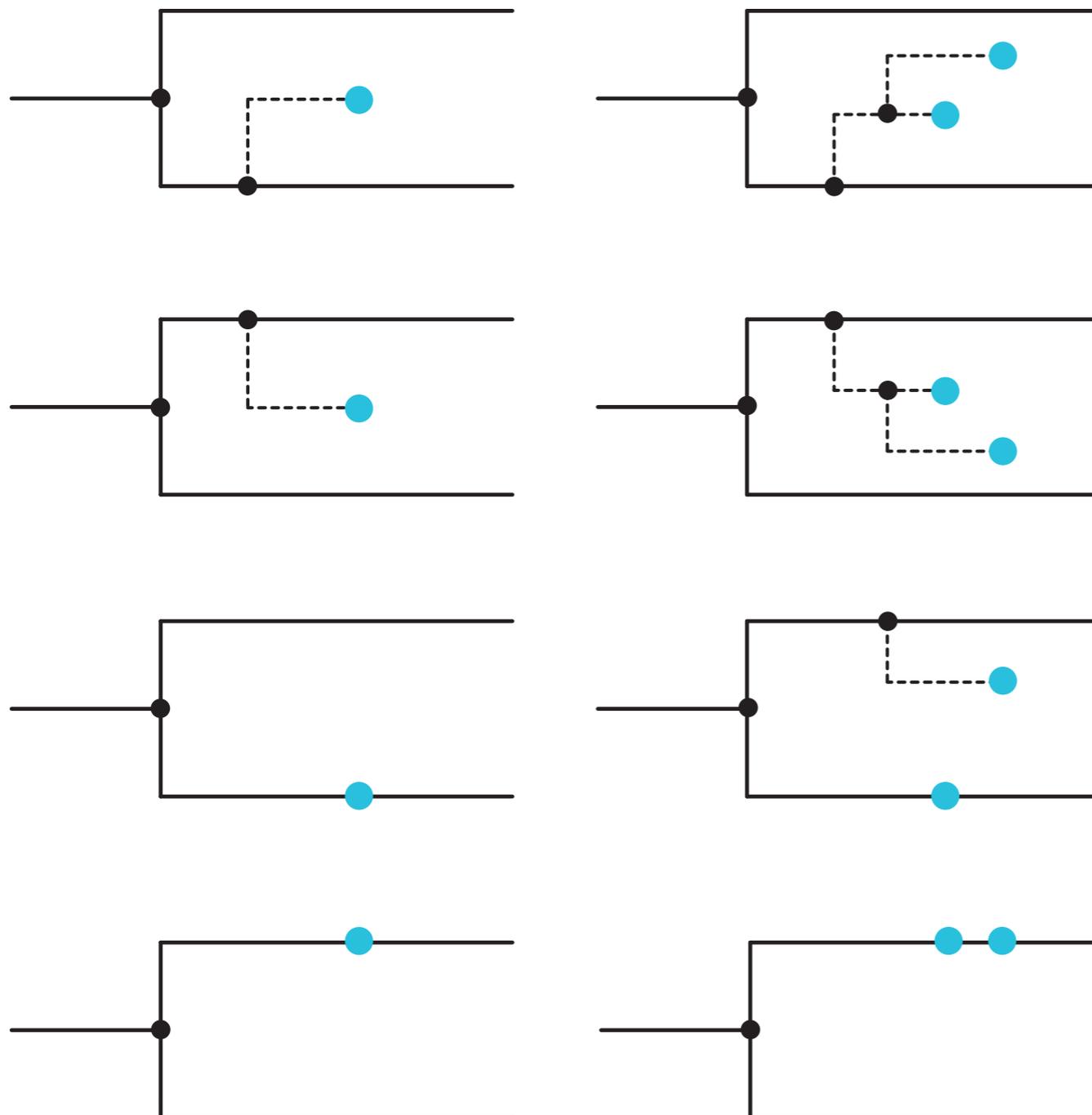
The fossilised birth-death process: taxonomic uncertainty



- MCMC is used to propose possible fossil placements
- rjMCMC is used to propose sampled ancestor placements

● Fossil occurrence
● Speciation event

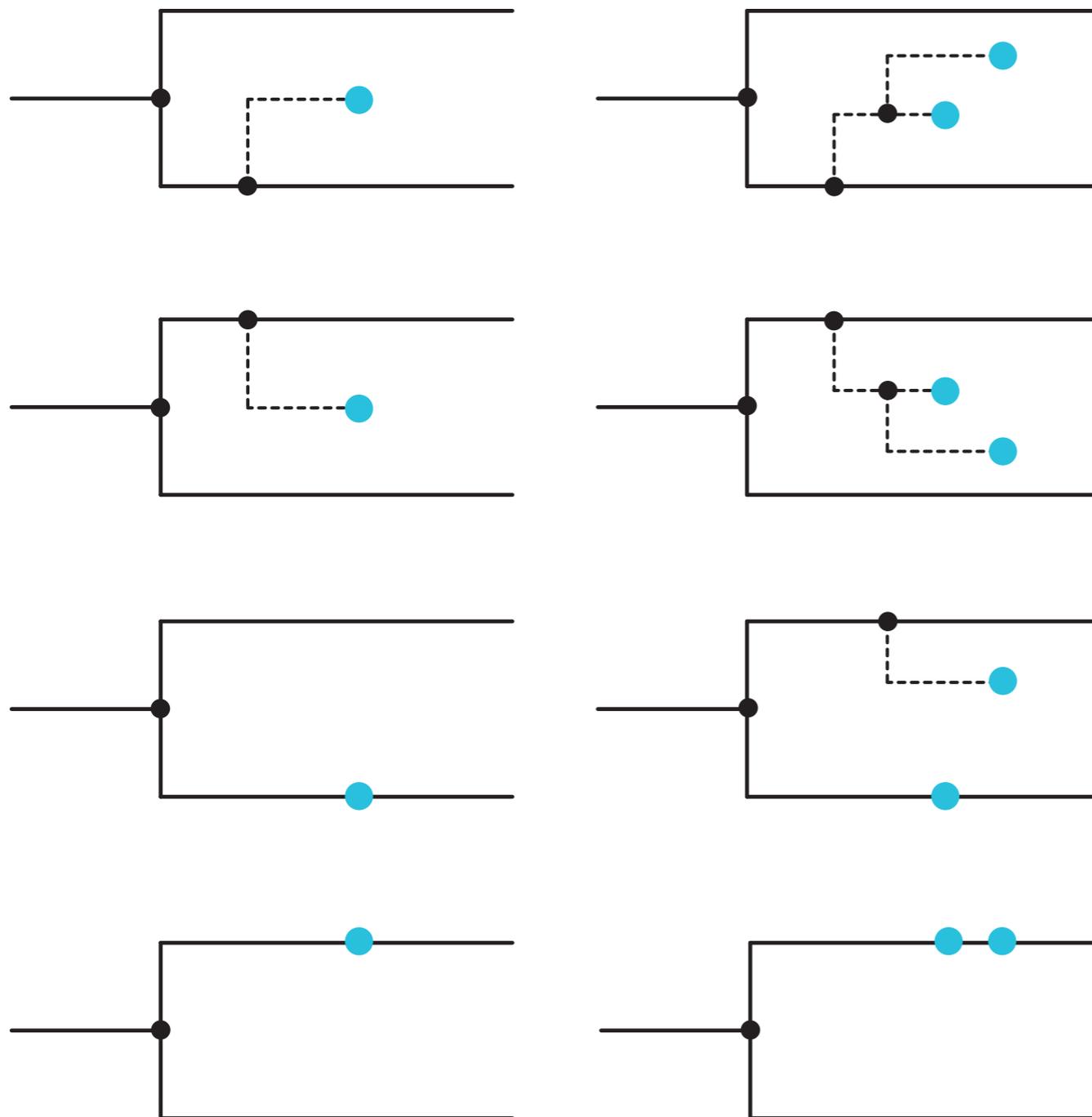
The fossilised birth-death process: taxonomic uncertainty



- Each fossil can attach anywhere along the tree, including along unobserved branches

● Fossil occurrence
● Speciation event

The fossilised birth-death process: taxonomic uncertainty



- The probability of any given realisation of the FBD process is conditional on the model parameters: λ , μ , ρ , ψ

● Fossil occurrence
● Speciation event

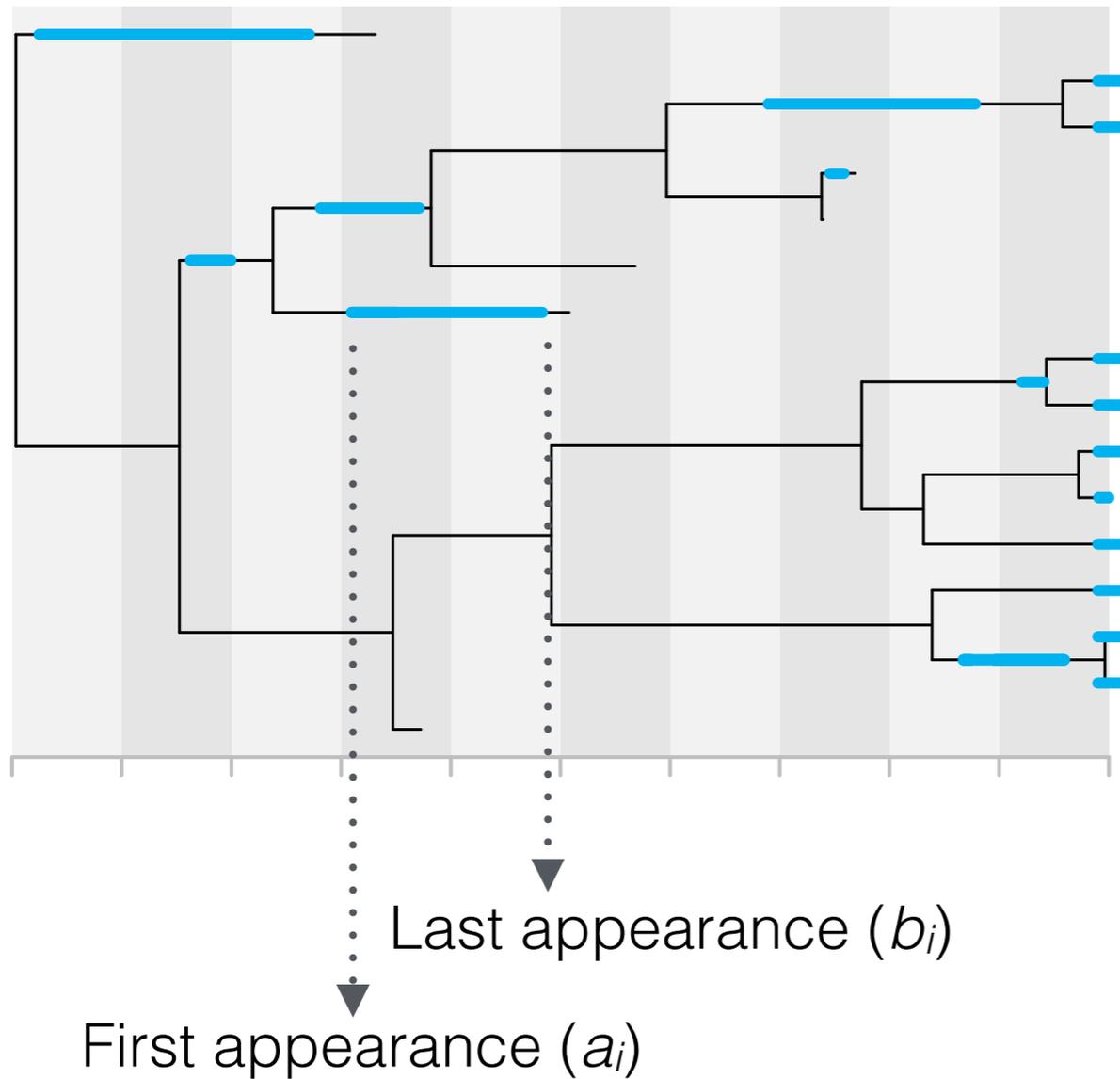
Application with or without character data for fossils

- Three alternative scenarios to applying to FBD model in divergence time estimation
 - Molecular data for extant taxa only ACAC...
TCAC...
ACAG...
 - Molecular data for extant taxa + morphological data for both extant & extinct taxa ACAC... 0101...
TCAC... 1101...
ACAG... 0100...
 - Morphological data for both extant and extinct taxa or extinct taxa only 0101...
1101...
0100...

Application with or without character data for fossils

- Three alternative scenarios to applying to FBD model in divergence time estimation
 - Molecular data for extant taxa only
ACAC...
TCAC...
ACAG...
 - Molecular data for extant taxa + morphological data for both extant & extinct taxa
ACAC... 0101...
TCAC... 1101...
ACAG... 0100...
 - Morphological data for both extant and extinct taxa or extinct taxa only
0101...
1101...
0100...
- The model can also be applied to estimate macroevolutionary parameters when you have no character data

Stratigraphic range data



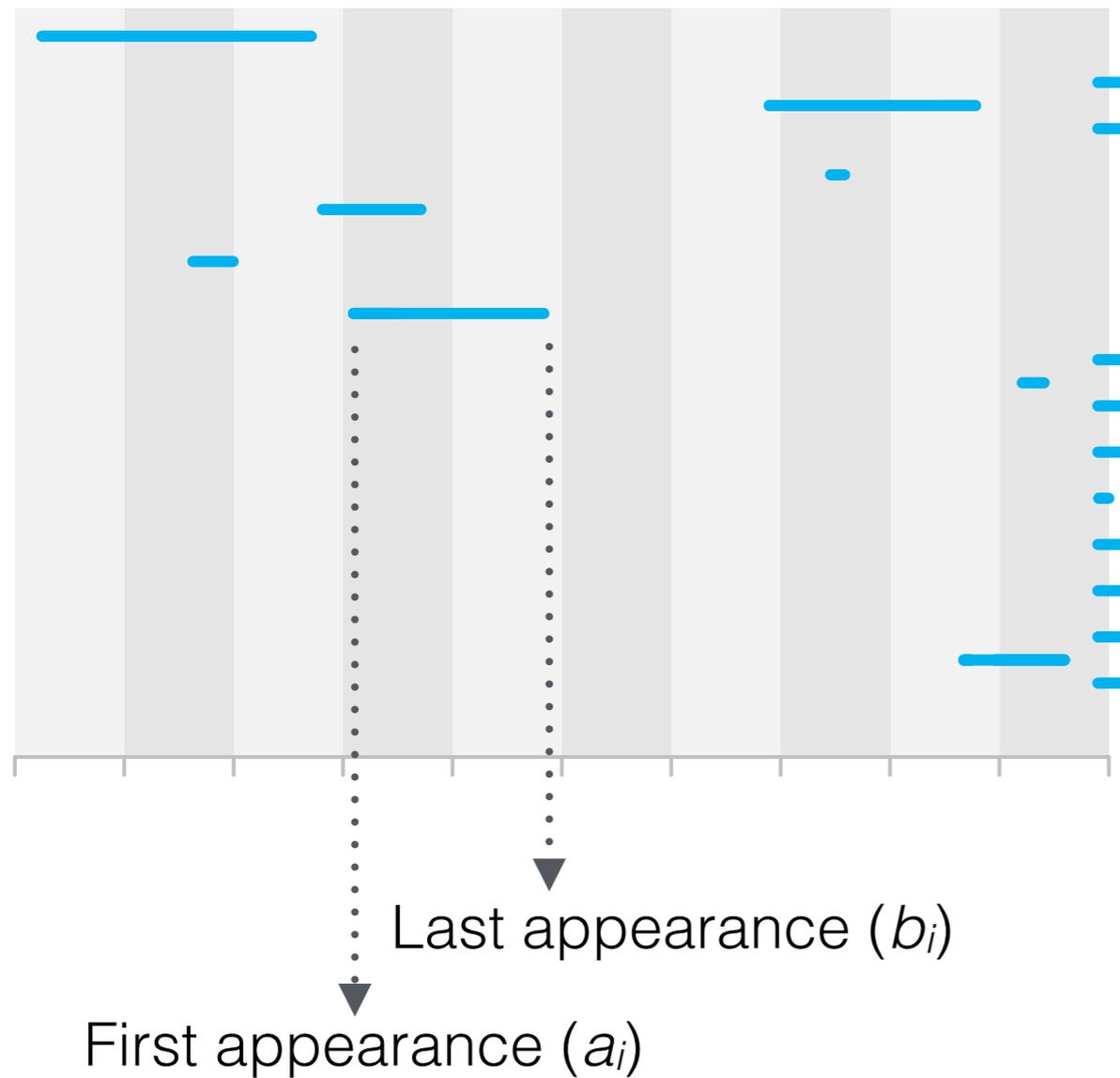
λ — speciation rate

μ — extinction rate

ψ — fossil sampling rate

ρ — extant species sampling

Stratigraphic range data



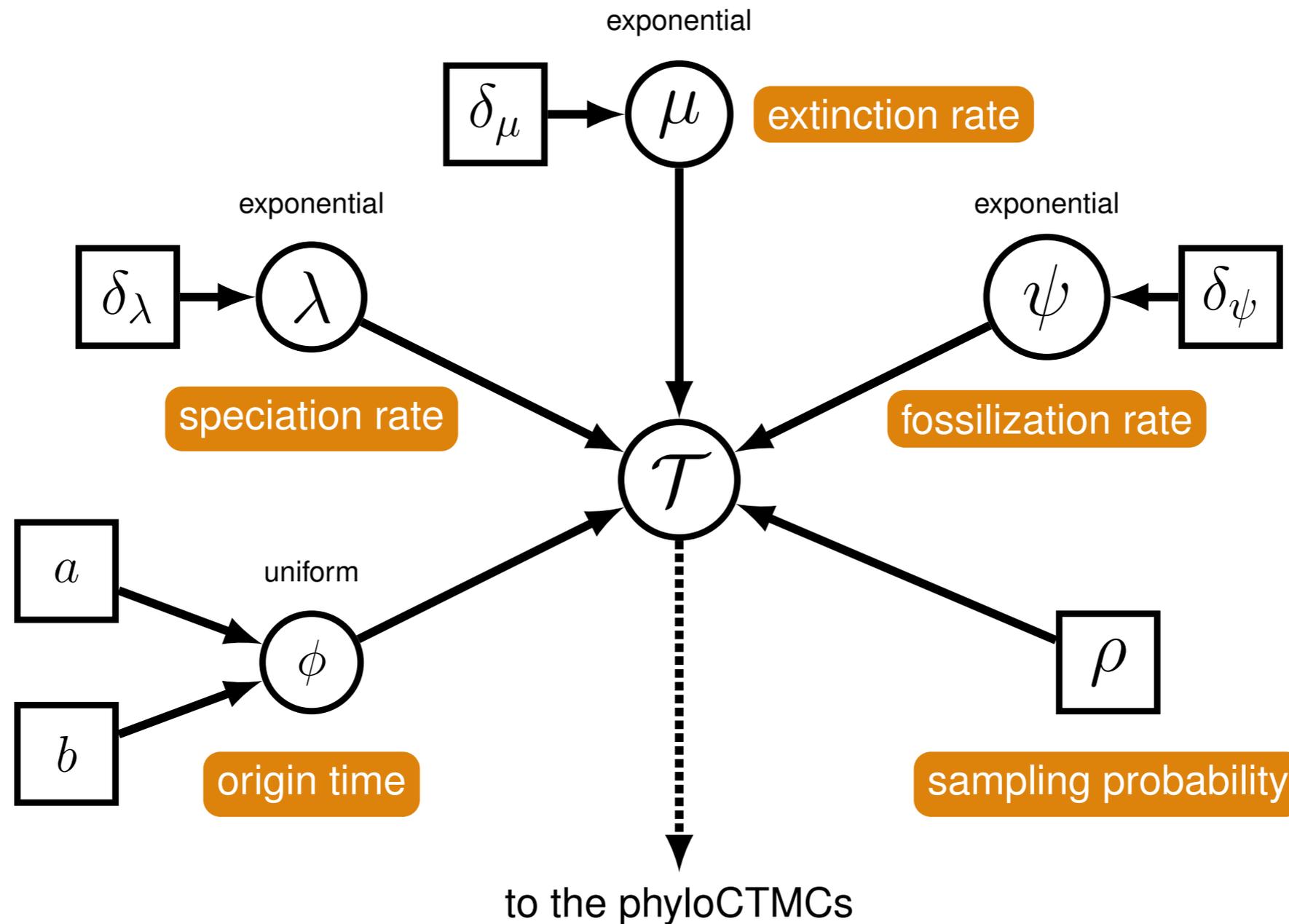
λ — speciation rate

μ — extinction rate

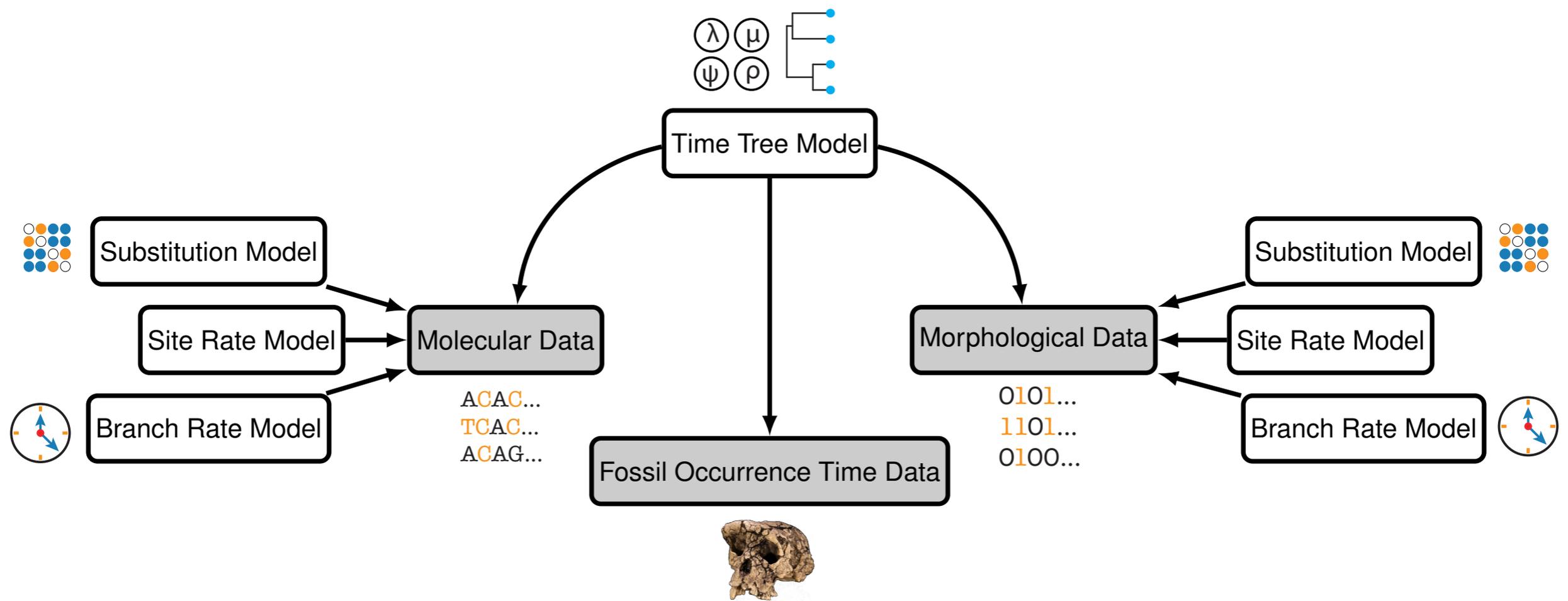
ψ — fossil sampling rate

ρ — extant species sampling

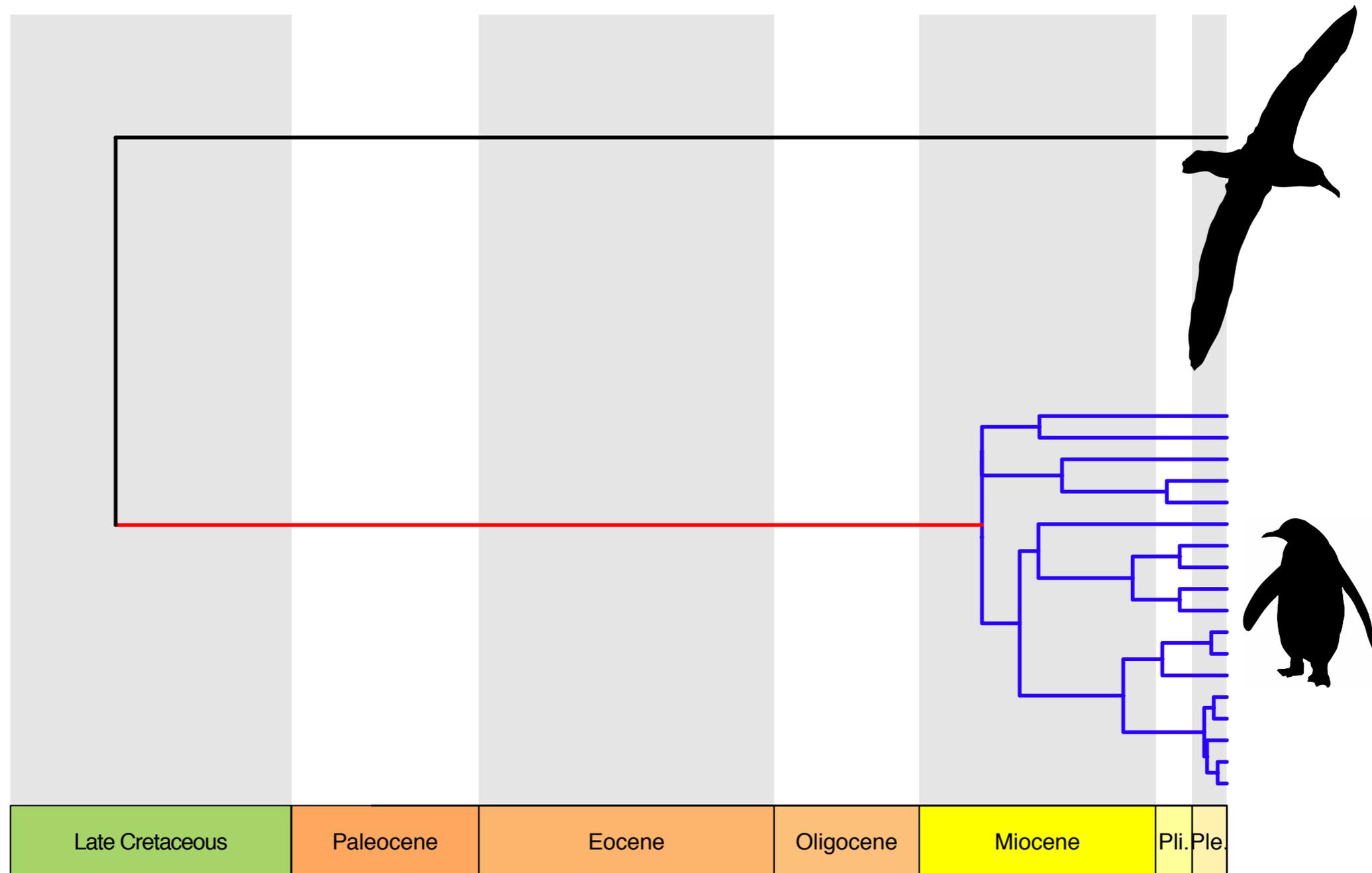
Putting things in a graphical modelling framework



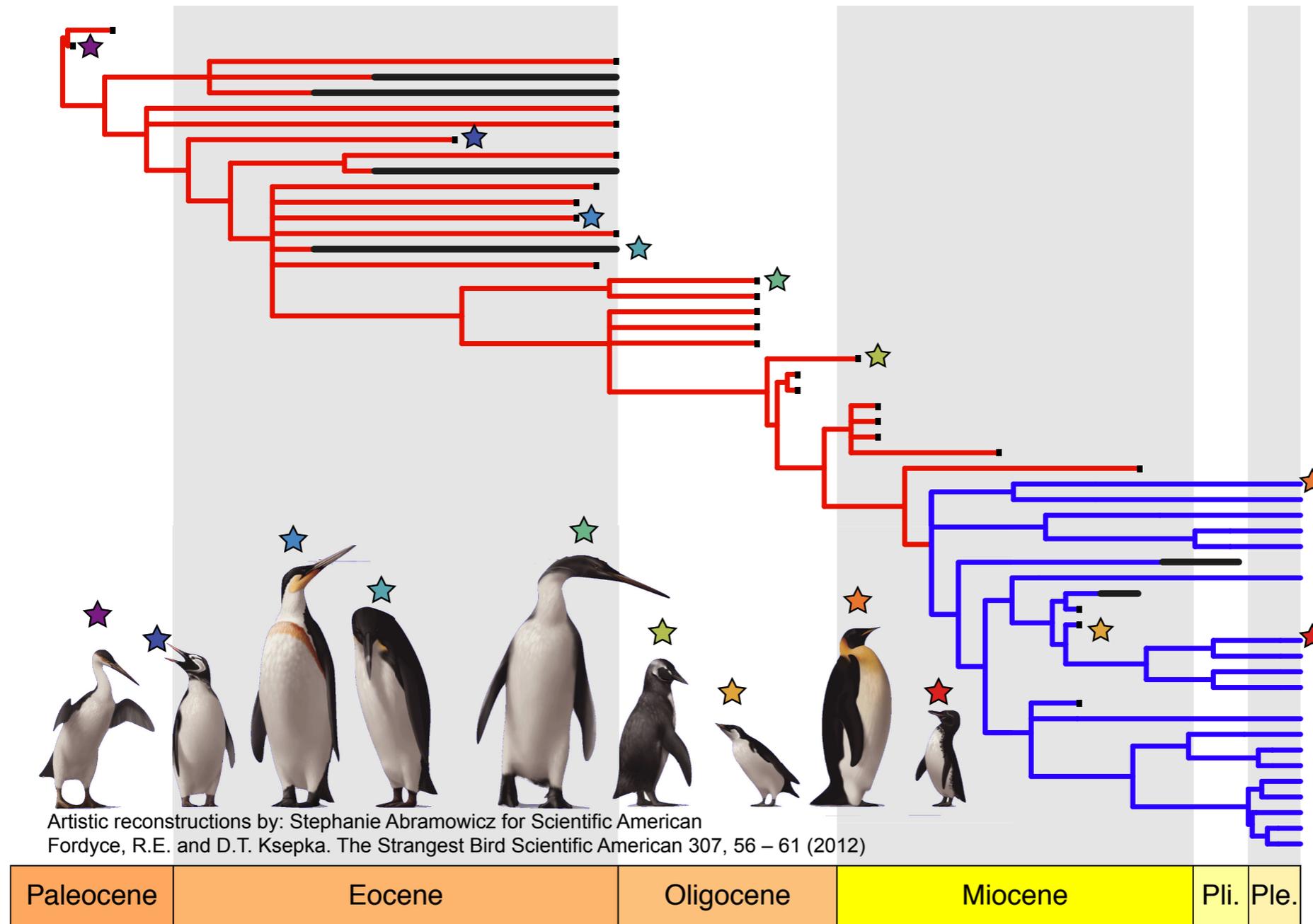
Putting things in a graphical modelling framework



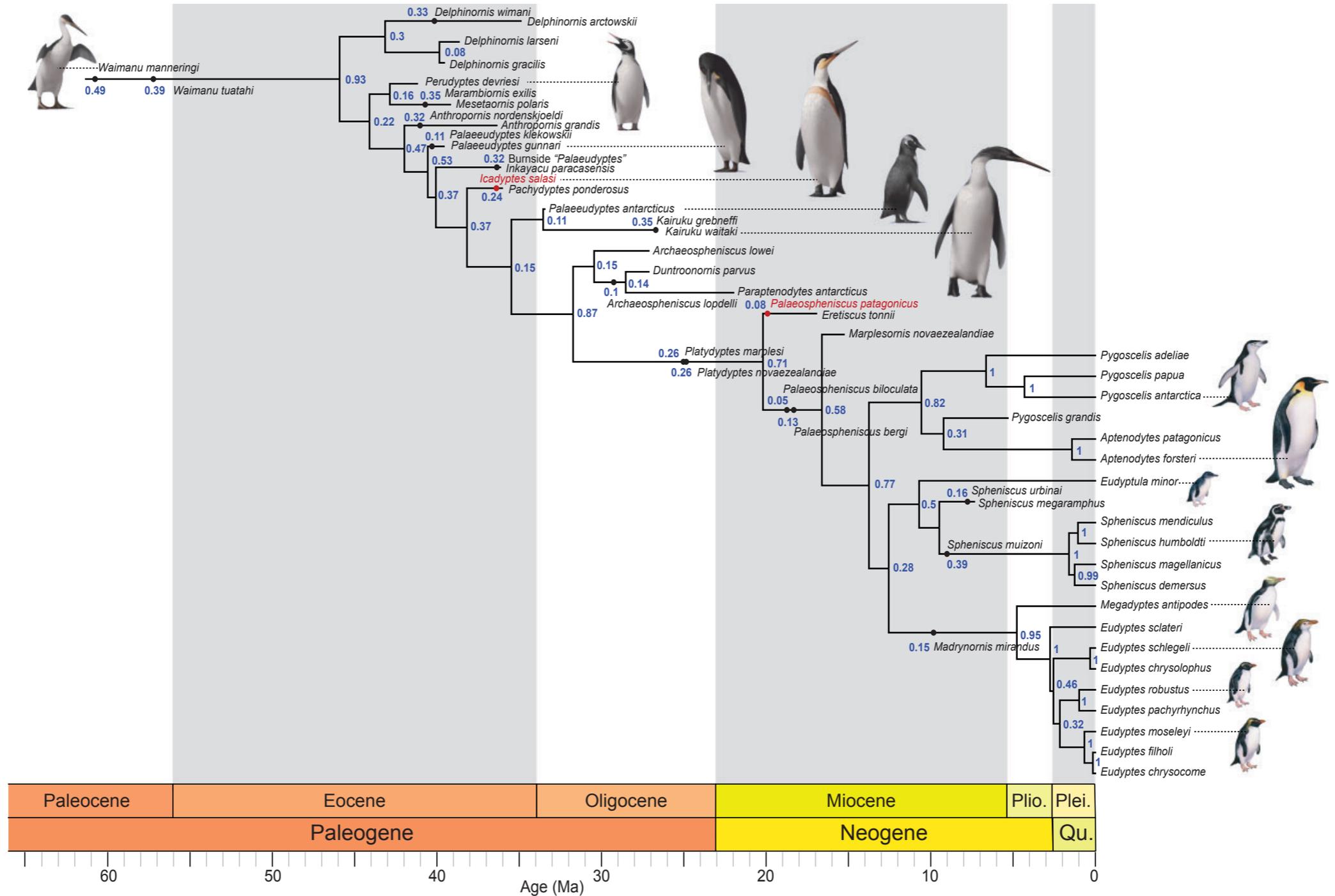
Dating the origin of crown penguins



Dating the origin of crown penguins



Dating the origin of crown penguins



Dating the origin of crown penguins

