

# BAYESIAN DIVERGENCE TIME ESTIMATION

Tracy Heath

Ecology, Evolution, & Organismal Biology  
Iowa State University

@tracy7  
<http://phyloworks.org>

2015 Workshop on Molecular Evolution  
Woods Hole, MA USA

# OUTLINE

## Overview of divergence time estimation

- Relaxed clock models – accounting for variation in substitution rates among lineages
- Tree models – lineage diversification and sampling

break

## BEAST v2.2.0 Tutorial — Divergence-time estimation under birth-death processes

<http://phyloworks.org/workshops/divtime.html>

### **Choose one:**

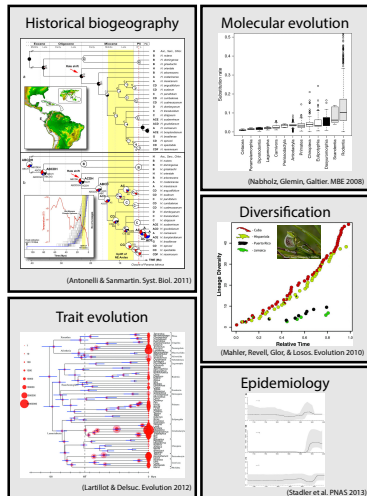
- Dating Bear Divergence Times with the Fossilized Birth-Death Process
- Estimating Epidemiological Parameters of an Ebola Outbreak

lobstah!

# A TIME-SCALE FOR EVOLUTION

## Phylogenetic divergence-time estimation

- What was the spacial and climatic environment of ancient angiosperms?
- How has mammalian body-size changed over time?
- How has the infection rate of HCV in Egypt changed over time?
- Is diversification in Caribbean anoles correlated with ecological opportunity?
- How has the rate of molecular evolution changed across the Tree of Life?



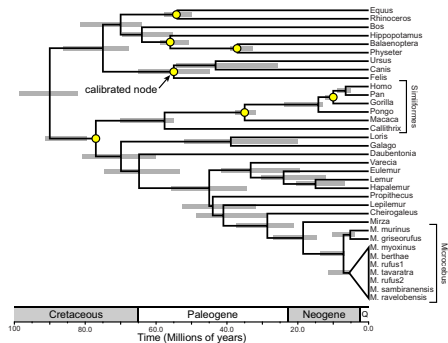
# DIVERGENCE TIME ESTIMATION

**Goal:** Estimate the ages of interior nodes to understand the timing and rates of evolutionary processes

Model how rates are distributed across the tree

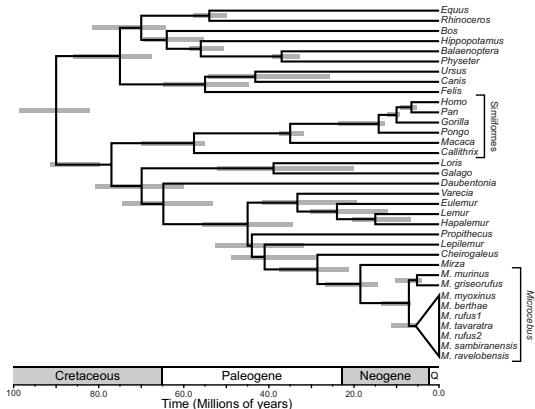
Describe the distribution of speciation events over time

External calibration information for estimates of absolute node times



# A TIME-SCALE FOR EVOLUTION

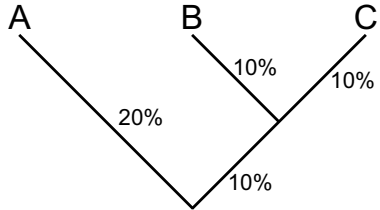
Phylogenetic trees can provide both topological information and temporal information



# THE GLOBAL MOLECULAR CLOCK

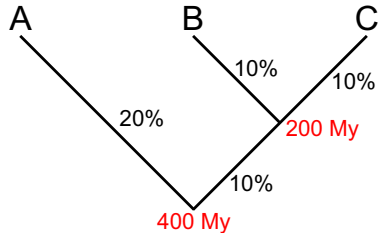
Assume that the rate of evolutionary change is constant over time

(branch lengths equal percent sequence divergence)



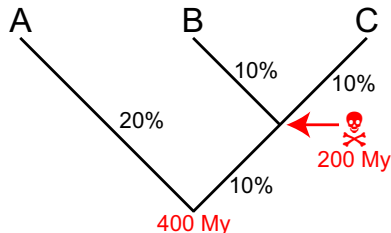
# THE GLOBAL MOLECULAR CLOCK

We can date the tree if we know the rate of change is 1% divergence per 10 My



# THE GLOBAL MOLECULAR CLOCK

If we found a fossil of the MRCA of **B** and **C**, we can use it to calculate the rate of change & date the root of the tree





# REJECTING THE GLOBAL MOLECULAR CLOCK

Rates of evolution vary across lineages and over time

## **Mutation rate:**

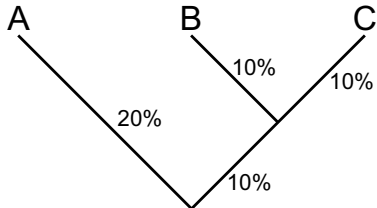
Variation in

- metabolic rate
- generation time
- DNA repair

## **Fixation rate:**

Variation in

- strength and targets of selection
- population sizes

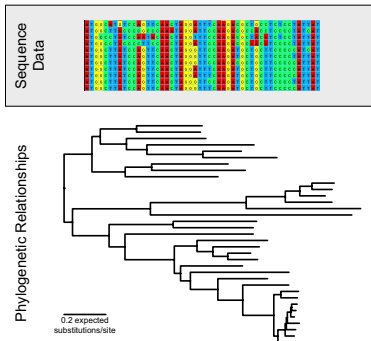


# UNCONSTRAINED ANALYSIS

Sequence data provide information about **branch lengths**

In units of **the expected # of substitutions per site**

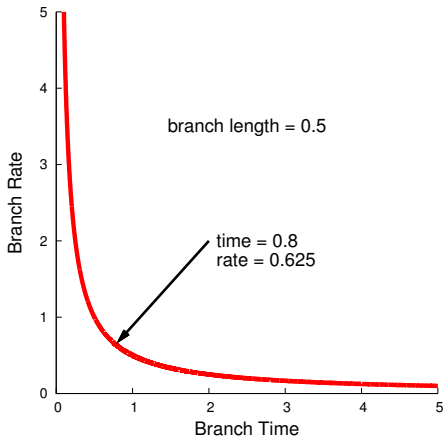
branch length = rate  $\times$  time



# RATE AND TIME

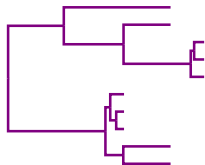
The sequence data  
provide information  
about branch length

for any possible rate,  
there's a time that fits  
the branch length  
perfectly

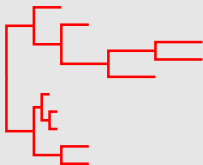


# RATE AND TIME

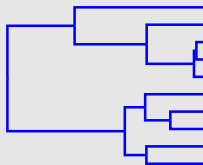
The expected # of substitutions/site occurring along a branch is the product of the substitution rate and time



length = rate  $\times$  time



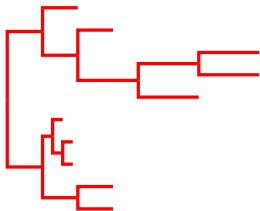
length = rate



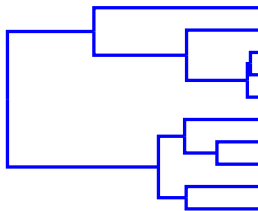
length = time

Methods for dating species divergences estimate the substitution rate and time separately

# BAYESIAN DIVERGENCE TIME ESTIMATION



length = rate



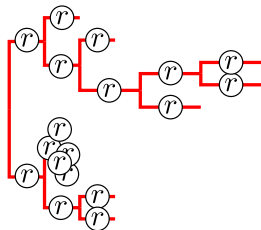
length = time

$$\mathcal{R} = (r_1, r_2, r_3, \dots, r_{2N-2})$$

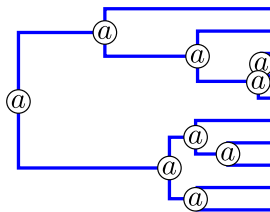
$$\mathcal{A} = (a_1, a_2, a_3, \dots, a_{N-1})$$

$$N = \text{number of tips}$$

# BAYESIAN DIVERGENCE TIME ESTIMATION



length = rate



length = time

$$\mathcal{R} = (r_1, r_2, r_3, \dots, r_{2N-2})$$

$$\mathcal{A} = (a_1, a_2, a_3, \dots, a_{N-1})$$

$$N = \text{number of tips}$$

# BAYESIAN DIVERGENCE TIME ESTIMATION

Posterior probability

$$f(\mathcal{R}, \mathcal{A}, \theta_{\mathcal{R}}, \theta_{\mathcal{A}}, \theta_s \mid D, \psi)$$

$\mathcal{R}$  Vector of rates on branches

$\mathcal{A}$  Vector of internal node ages

$\theta_{\mathcal{R}}, \theta_{\mathcal{A}}, \theta_s$  Model parameters

$D$  Sequence data

$\psi$  Tree topology

# BAYESIAN DIVERGENCE TIME ESTIMATION

$$f(\mathcal{R}, \mathcal{A}, \theta_{\mathcal{R}}, \theta_{\mathcal{A}}, \theta_s | D) =$$

$$\frac{f(D | \mathcal{R}, \mathcal{A}, \theta_s) f(\mathcal{R} | \theta_{\mathcal{R}}) f(\mathcal{A} | \theta_{\mathcal{A}}) f(\theta_s)}{f(D)}$$

$$f(D | \mathcal{R}, \mathcal{A}, \theta_{\mathcal{R}}, \theta_{\mathcal{A}}, \theta_s)$$

Likelihood

$$f(\mathcal{R} | \theta_{\mathcal{R}})$$

Prior on rates

$$f(\mathcal{A} | \theta_{\mathcal{A}})$$

Prior on node ages

$$f(\theta_s)$$

Prior on substitution parameters

$$f(D)$$

Marginal probability of the data



# BAYESIAN DIVERGENCE TIME ESTIMATION

Estimating divergence times relies on 2 main elements:

- Branch-specific rates:  $f(\mathcal{R} \mid \theta_{\mathcal{R}})$
- Node ages:  $f(\mathcal{A} \mid \theta_{\mathcal{A}}, \mathcal{C})$

# MODELING RATE VARIATION

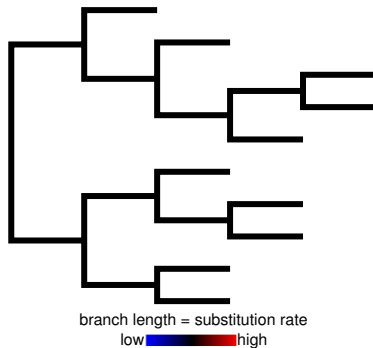
Some models describing lineage-specific substitution rate variation:

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

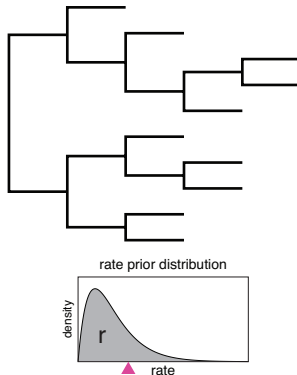
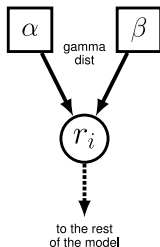
# GLOBAL MOLECULAR CLOCK

The substitution rate is constant over time

All lineages share the same rate



# GLOBAL MOLECULAR CLOCK



# RELAXED-CLOCK MODELS

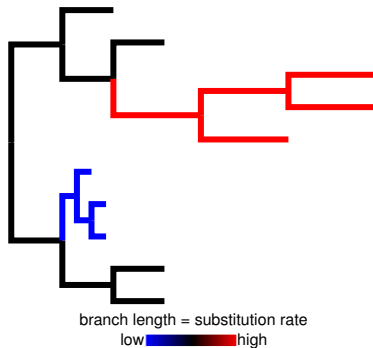
To accommodate variation in substitution rates  
'relaxed-clock' models estimate lineage-specific substitution rates

- **Local molecular clocks**
- **Punctuated rate change model**
- **Log-normally distributed autocorrelated rates**
- **Uncorrelated/independent rates models**
- **Mixture models on branch rates**

# LOCAL MOLECULAR CLOCKS

Rate shifts occur  
infrequently over the tree

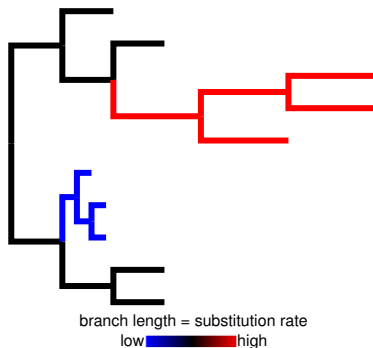
Closely related lineages  
have equivalent rates  
(clustered by sub-clades)



# LOCAL MOLECULAR CLOCKS

Most methods for estimating local clocks required specifying the number and locations of rate changes *a priori*

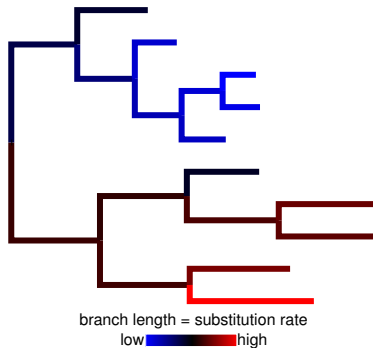
Drummond and Suchard (2010) introduced a Bayesian method that samples over a broad range of possible *random local clocks*



# AUTOCORRELATED RATES

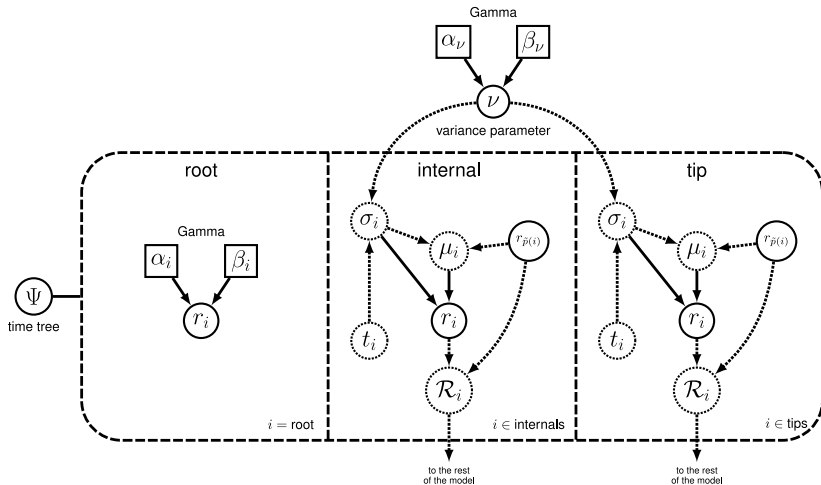
Substitution rates evolve gradually over time — closely related lineages have similar rates

The rate at a node is drawn from a lognormal distribution with a mean equal to the parent rate





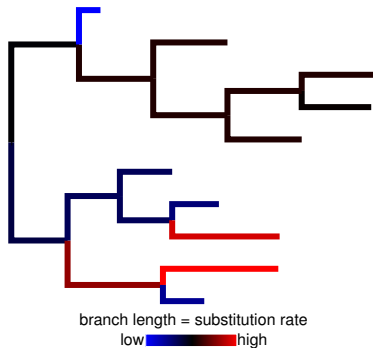
# AUTOCORRELATED RATES



# PUNCTUATED RATE CHANGE

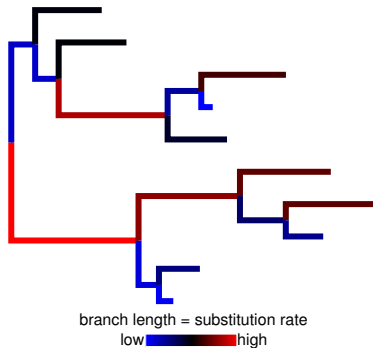
Rate changes occur along lineages according to a point process

At rate-change events, the new rate is a product of the parent's rate and a  $\Gamma$ -distributed multiplier



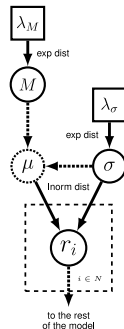
# INDEPENDENT/UNCORRELATED RATES

Lineage-specific rates are uncorrelated when the rate assigned to each branch is independently drawn from an underlying distribution



# INDEPENDENT/UNCORRELATED RATES

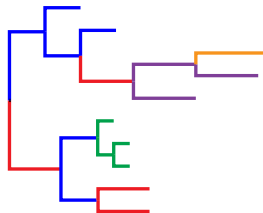
Lineage-specific rates are uncorrelated when the rate assigned to each branch is independently drawn from an underlying distribution



# INFINITE MIXTURE MODEL

## Dirichlet process prior:

Branches are partitioned into distinct rate categories



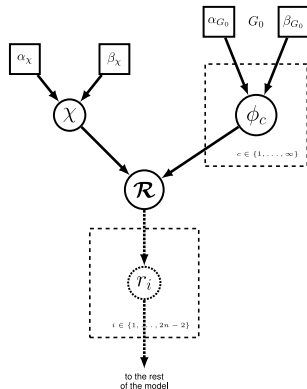
branch length = substitution rate



substitution rate classes

# INFINITE MIXTURE MODEL

**Dirichlet process prior:**  
Branches are partitioned  
into distinct rate categories



# MODELING RATE VARIATION

These are only a subset of the available models for branch-rate variation

- **Global molecular clock**
- **Local molecular clocks**
- **Punctuated rate change model**
- **Log-normally distributed autocorrelated rates**
- **Uncorrelated/independent rates models**
- **Dirchlet process prior**





# MODELING RATE VARIATION

These are only a subset of the available models for branch-rate variation

- **Global molecular clock**
- **Local molecular clocks**
- **Punctuated rate change model**
- **Log-normally distributed autocorrelated rates**
- **Uncorrelated/independent rates models**
- **Dirichlet process prior**

Considering model selection, uncertainty, & plausibility is **very** important for Bayesian divergence time analysis



# BAYESIAN DIVERGENCE TIME ESTIMATION

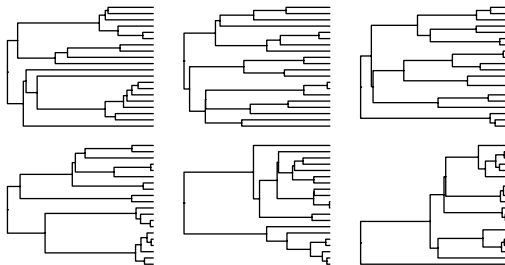
Estimating divergence times relies on 2 main elements:

- Branch-specific rates:  $f(\mathcal{R} \mid \theta_{\mathcal{R}})$
- Node ages:  $f(\mathcal{A} \mid \theta_{\mathcal{A}})$

<http://bayesiancook.blogspot.com/2013/12/two-sides-of-same-coin.html>

# PRIORS ON THE TREE AND NODE AGES

Relaxed clock Bayesian analyses require a prior distribution on time trees



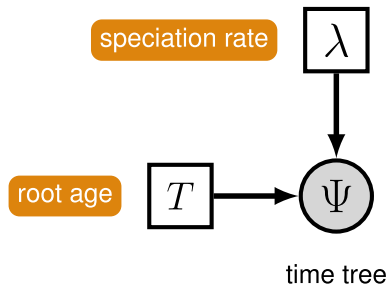
Different node-age priors make different assumptions about the timing of divergence events

# STOCHASTIC BRANCHING PROCESSES

Node-age priors based on stochastic models of lineage diversification

**Yule process:** assumes a constant rate of speciation, across lineages

A pure birth process—every node leaves extant descendants (no extinction)

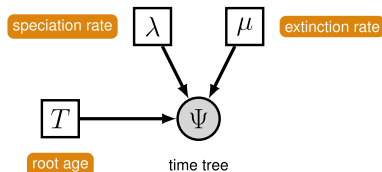


# STOCHASTIC BRANCHING PROCESSES

Node-age priors based on stochastic models of lineage diversification

## Constant-rate birth-death

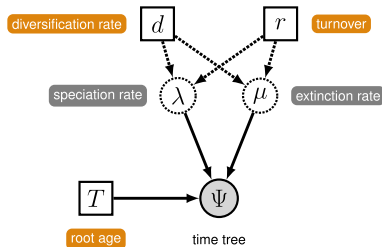
**process:** at any point in time a lineage can speciate at rate  $\lambda$  or go extinct with a rate of  $\mu$



# STOCHASTIC BRANCHING PROCESSES

Node-age priors based on stochastic models of lineage diversification

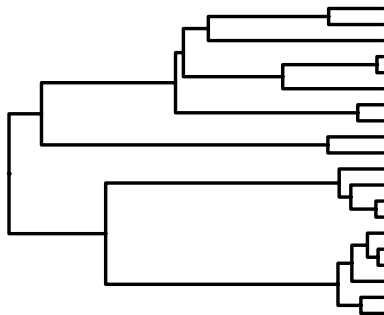
**Constant-rate birth-death process:** at any point in time a lineage can speciate at rate  $\lambda$  or go extinct with a rate of  $\mu$



# STOCHASTIC BRANCHING PROCESSES

Node-age priors based on stochastic models of lineage diversification

**Constant-rate birth-death process:** at any point in time a lineage can speciate at rate  $\lambda$  or go extinct with a rate of  $\mu$

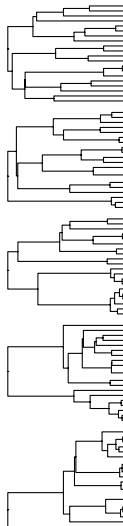


# STOCHASTIC BRANCHING PROCESSES

Different values of  $\lambda$  and  $\mu$  lead to different trees

Bayesian inference under these models can be very sensitive to the values of these parameters

Using hyperpriors on  $\lambda$  and  $\mu$  accounts for uncertainty in these hyperparameters





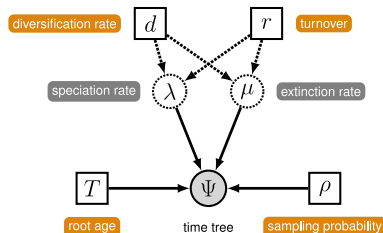
# STOCHASTIC BRANCHING PROCESSES

Node-age priors based on stochastic models of lineage diversification

## Birth-death-sampling

**process:** an extension of the constant-rate birth-death model that accounts for random sampling of tips

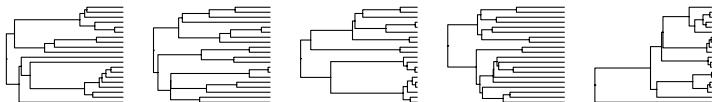
Conditions on a probability of sampling a tip,  $\rho$



# PRIORS ON NODE TIMES

Sequence data are only informative on *relative* rates & times

Node-time priors cannot give precise estimates of *absolute* node ages



We need external information (like fossils) to provide absolute time scale

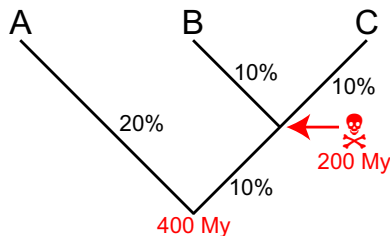


# CALIBRATING DIVERGENCE TIMES

Fossils (or other data) are necessary to estimate *absolute* node ages

There is **no information** in the sequence data for absolute time

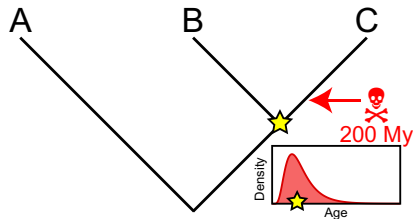
Uncertainty in the placement of fossils



# CALIBRATION DENSITIES

Bayesian inference is well suited to accommodating uncertainty in the age of the calibration node

Divergence times are calibrated by placing parametric densities on internal nodes offset by age estimates from the fossil record

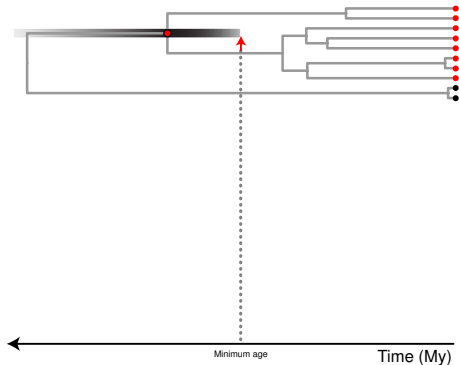




# FOSSIL CALIBRATION

Age estimates from fossils can provide **minimum** time constraints for internal nodes

Reliable **maximum** bounds are typically unavailable

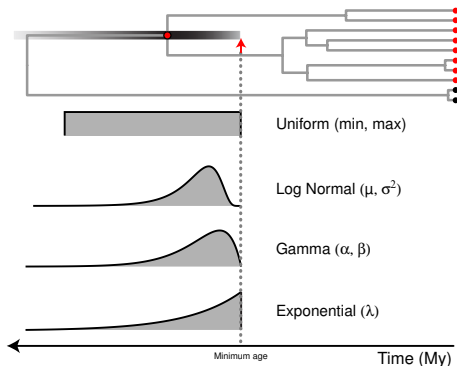


# PRIOR DENSITIES ON CALIBRATED NODES

## Common practice in Bayesian divergence-time estimation:

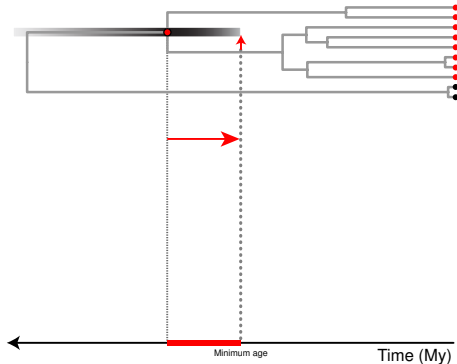
Parametric distributions are typically off-set by the age of the oldest fossil assigned to a clade

These prior densities do not (necessarily) require specification of maximum bounds



# PRIOR DENSITIES ON CALIBRATED NODES

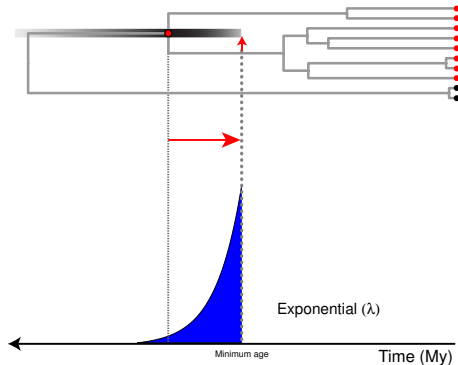
Describe the waiting time between the divergence event and the age of the oldest fossil





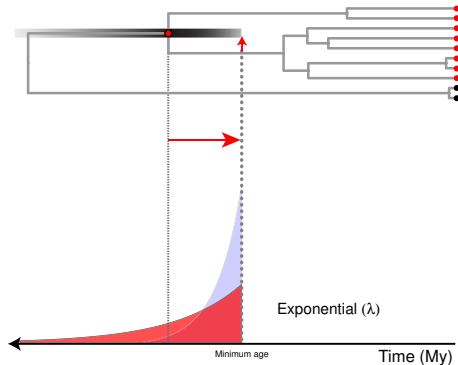
# PRIOR DENSITIES ON CALIBRATED NODES

Overly **informative** priors can bias node age estimates to be too young



# PRIOR DENSITIES ON CALIBRATED NODES

Uncertainty in the age of the MRCA of the clade relative to the age of the fossil may be better captured by **vague** prior densities

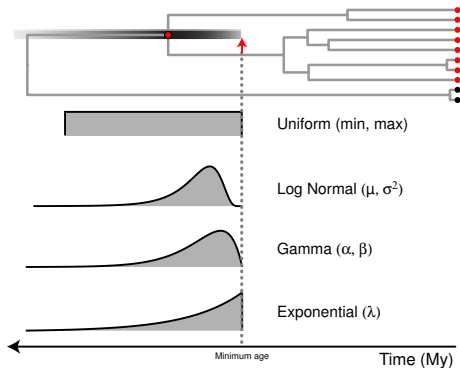


# PRIOR DENSITIES ON CALIBRATED NODES

## Common practice in Bayesian divergence-time estimation:

Estimates of absolute node ages are driven primarily by the calibration density

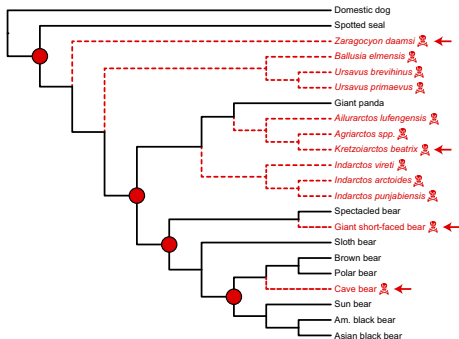
Specifying appropriate densities is a challenge for most molecular biologists



# IMPROVING FOSSIL CALIBRATION

We would prefer to eliminate the need for *ad hoc* calibration prior densities

Calibration densities do not account for diversification of fossils

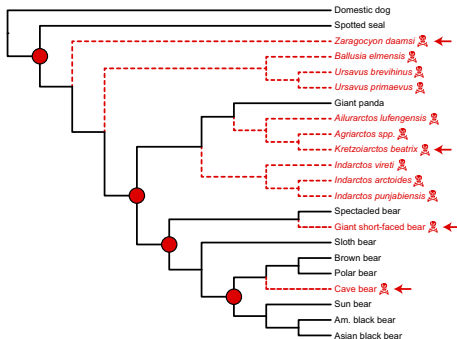


# IMPROVING FOSSIL CALIBRATION

We want to use all  
of the available fossils

## Example: Bears

12 fossils are reduced  
to 4 calibration ages  
with calibration density  
methods

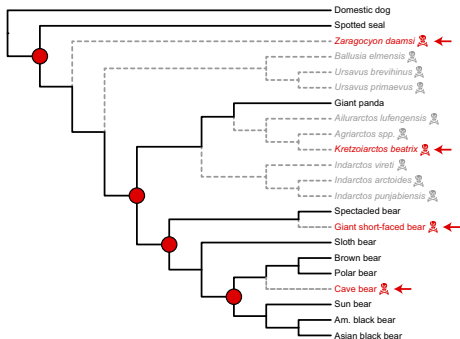


# IMPROVING FOSSIL CALIBRATION

We want to use all  
of the available fossils

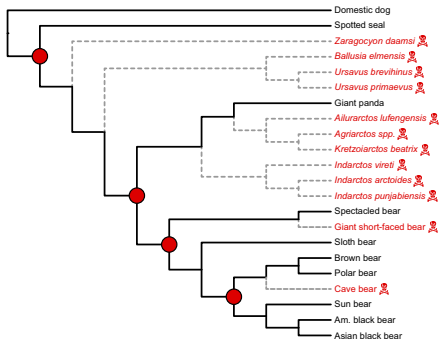
## Example: Bears

12 fossils are reduced  
to 4 calibration ages  
with calibration density  
methods



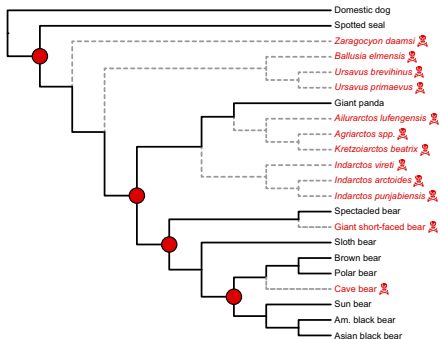
# IMPROVING FOSSIL CALIBRATION

Because fossils are part of the diversification process, we can combine fossil calibration with birth-death models



# IMPROVING FOSSIL CALIBRATION

This relies on a branching model that accounts for **speciation, extinction, and rates of fossilization, preservation, and recovery**





# THE FOSSILIZED BIRTH-DEATH PROCESS (FBD)

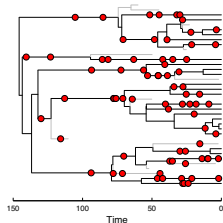
## Improving statistical inference of absolute node ages

Eliminates the need to specify arbitrary calibration densities

Better capture our statistical uncertainty in species divergence dates

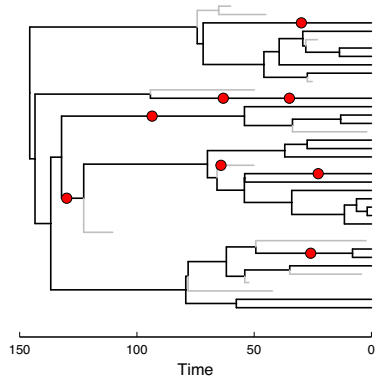
**All** reliable fossils associated with a clade are used

Useful for calibration or 'total-evidence' dating



# THE FOSSILIZED BIRTH-DEATH PROCESS (FBD)

Recovered fossil specimens provide historical observations of the diversification process that generated the tree of extant species



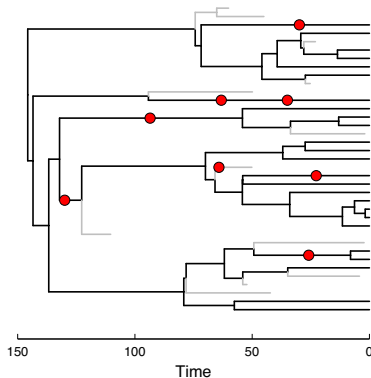
# THE FOSSILIZED BIRTH-DEATH PROCESS (FBD)

The probability of the tree and fossil observations under a birth-death model with rate parameters:

$\lambda$  = speciation

$\mu$  = extinction

$\psi$  = fossilization/recovery



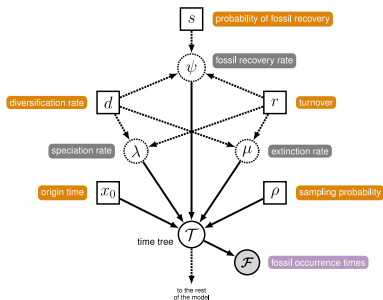
# THE FOSSILIZED BIRTH-DEATH PROCESS (FBD)

The probability of the tree and fossil observations under a birth-death model with rate parameters:

$\lambda$  = speciation

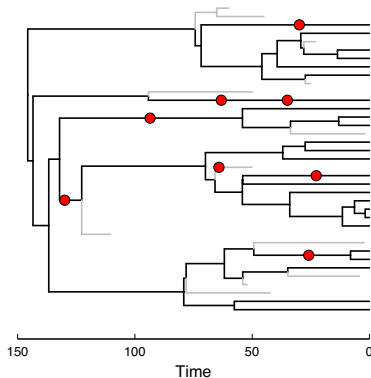
$\mu$  = extinction

$\psi$  = fossilization/recovery



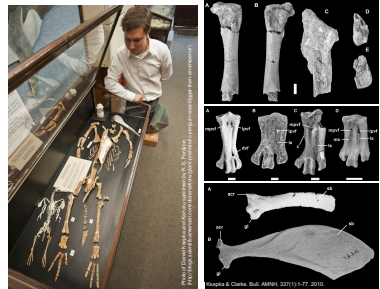
# THE FOSSILIZED BIRTH-DEATH PROCESS (FBD)

We use MCMC to sample realizations of the diversification process, integrating over the topology—including placement of the fossils—and speciation times

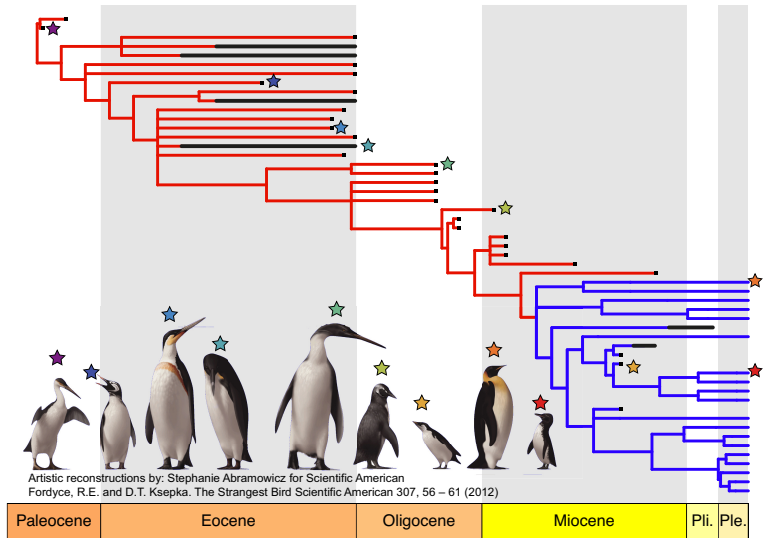


# PENGUIN DIVERSITY IN DEEP TIME

Can we improve our understanding of penguin evolution by considering both extant and fossil taxa?

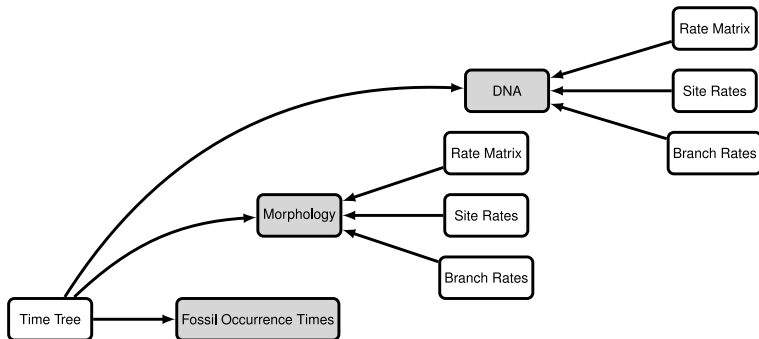


# PENGUIN DIVERSITY IN DEEP TIME



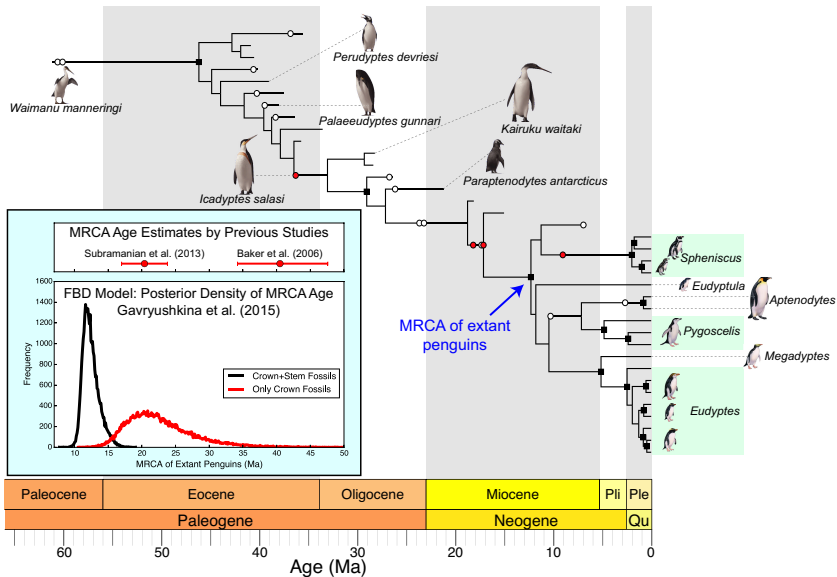
# INTEGRATIVE BAYESIAN INFERENCE

Combine models for DNA sequence evolution, morphological change, and fossil recovery over time to jointly estimate the tree topology, divergence times, and lineage diversification rates



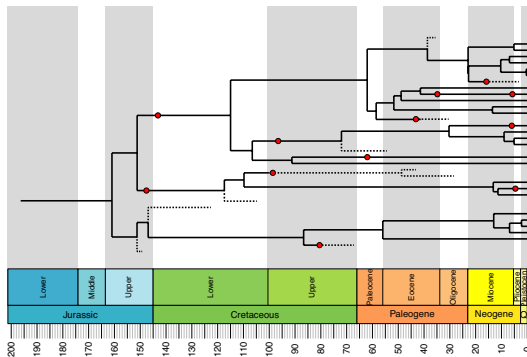


# PENGUIN DIVERSITY IN DEEP TIME



# INFERRING FBD TREES

Extensions of the fossilized birth-death process accommodate variation in fossil sampling, non-random species sampling, & shifts in diversification rates.

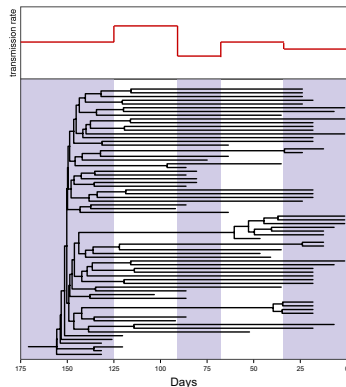


With character data for both fossil & extant species, we account for uncertainty in fossil placement

# SKYLINE BIRTH-DEATH PROCESS

A piecewise shifting model where parameters change over time

Used to estimate epidemiological parameters of an outbreak



**Birth–death skyline plot reveals temporal changes of epidemic spread in HIV and hepatitis C virus (HCV)**

Tanja Stadler<sup>a,1,2</sup>, Denise Kühnert<sup>b,c,1</sup>, Sebastian Bonhoeffer<sup>a</sup>, and Alexei J. Drummond<sup>b,c</sup>

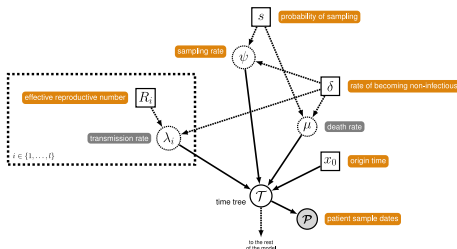
# SKYLINE BIRTH-DEATH PROCESS

$l$  is the number of parameter intervals

$R_i$  is the effective reproductive number for interval  $i \in l$

$\delta$  is the rate of becoming non-infectious

$s$  is the probability of sampling an individual after becoming non-infectious



$$R_i = \frac{\lambda_i}{\mu + \psi}, \quad \delta = \mu + \psi, \quad s = \frac{\psi}{\mu + \psi}$$

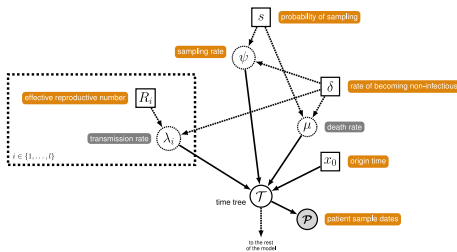
# SKYLINE BIRTH-DEATH PROCESS

$l$  is the number of parameter intervals

$\lambda_i$  is the transmission rate for interval  $i \in l$

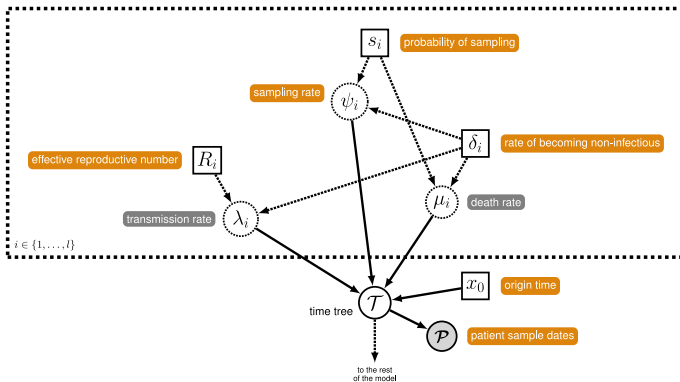
$\mu$  is the viral lineage death rate

$\psi$  is the rate each individual is sampled



$$\lambda_i = R_i \delta, \quad \mu = \delta - s\delta, \quad \psi = s\delta$$

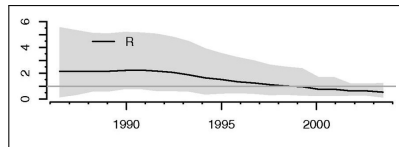
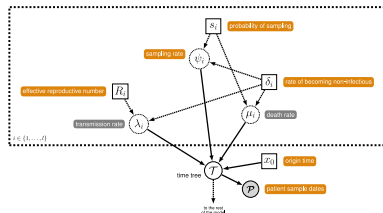
# SKYLINE BIRTH-DEATH PROCESS



# SKYLINE BIRTH-DEATH PROCESS

A decline in  $R$  over the history of HIV-1 in the UK is consistent with the introduction of effective drug therapies

After 1998  $R$  decreased below 1, indicating a declining epidemic



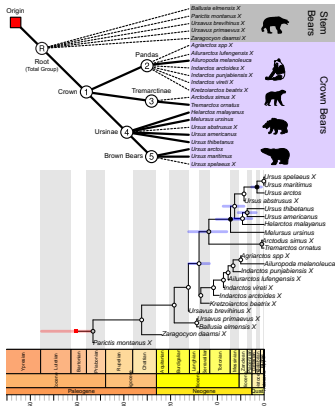
# QUESTIONS?





# EXERCISES: CHOOSE YOUR OWN ADVENTURE

## Dating Bear Divergence Times with the Fossilized Birth-Death Process



## Estimating Epidemiological Parameters of an Ebola Outbreak

