

Tip-dating and phylodynamics

Sebastian Duchene

In this talk...

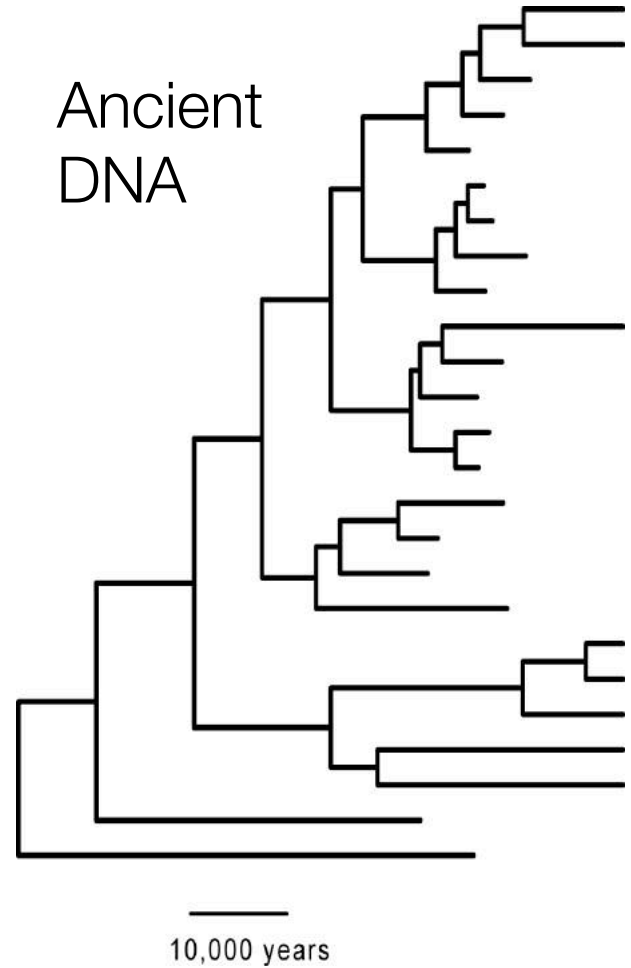
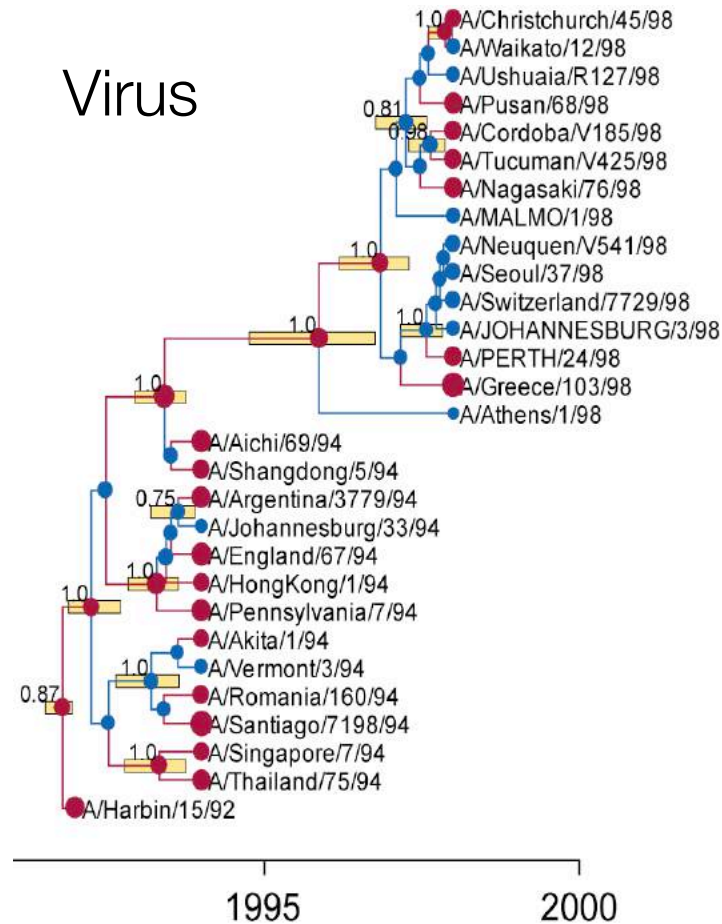
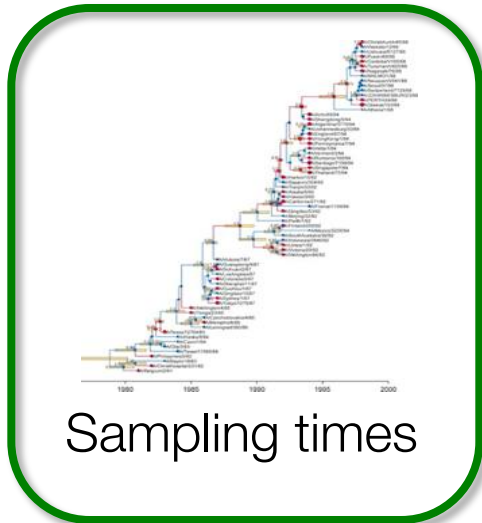
1. Measurably evolving populations
2. Assessing temporal structure
3. Bayesian methods to assess temporal structure
 1. Case studies in influenza viruses and *Mycobacterium tuberculosis*
4. A case study using ancient DNA
5. Tree priors and tip-dating
 1. Case studies in influenza and human respiratory syncytial virus
6. Model adequacy in phylodynamics
 1. Assessing the sampling process

Measurably Evolving Populations

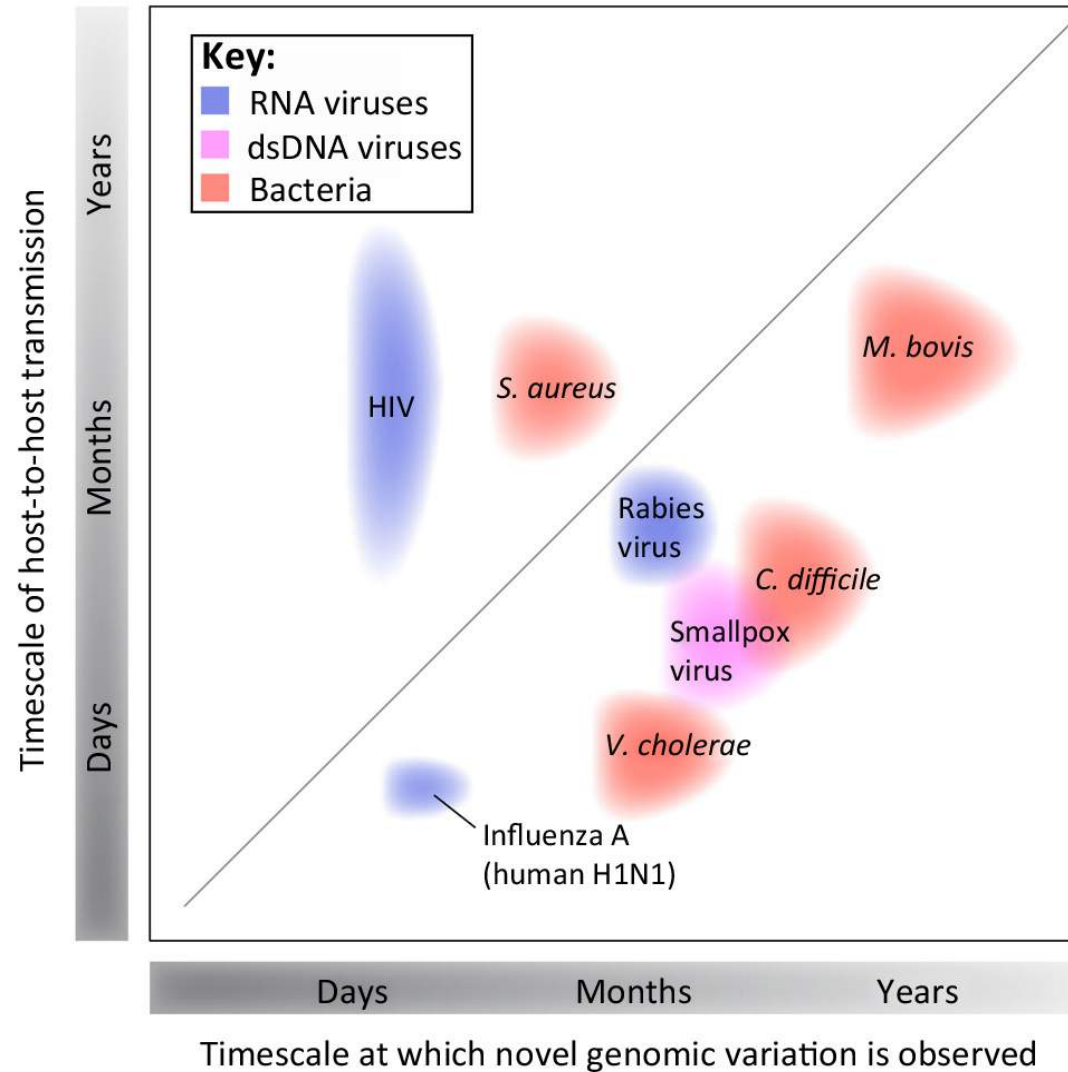
Measurably evolving populations

- Can perform tip-dating on measurably evolving populations
 - Substantial genetic change during the sampling window
 - High rate (pathogens) or wide window (ancient DNA)
- Sampling window represents large fraction of the tree height

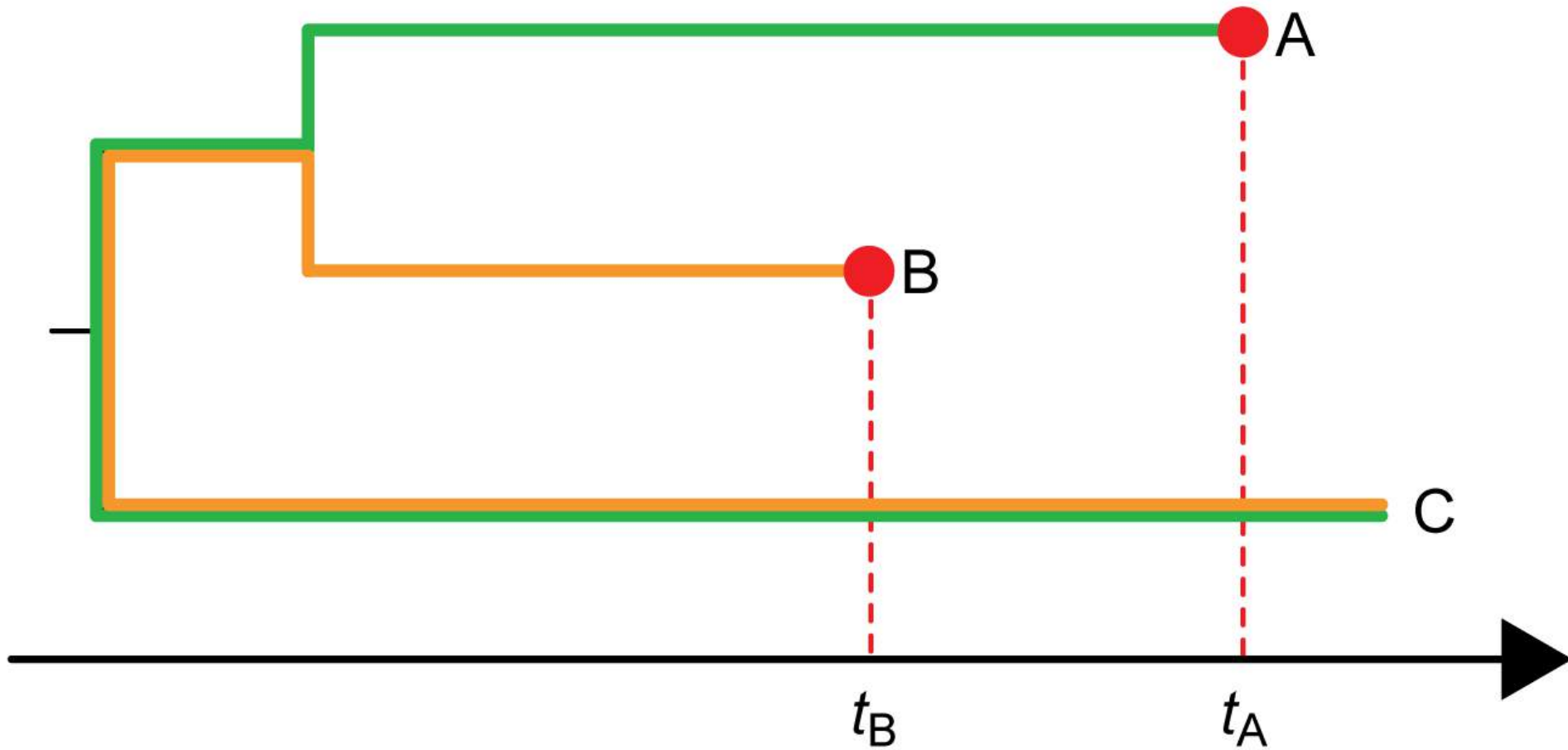
Calibration: Sampling times



Rapidly evolving pathogens



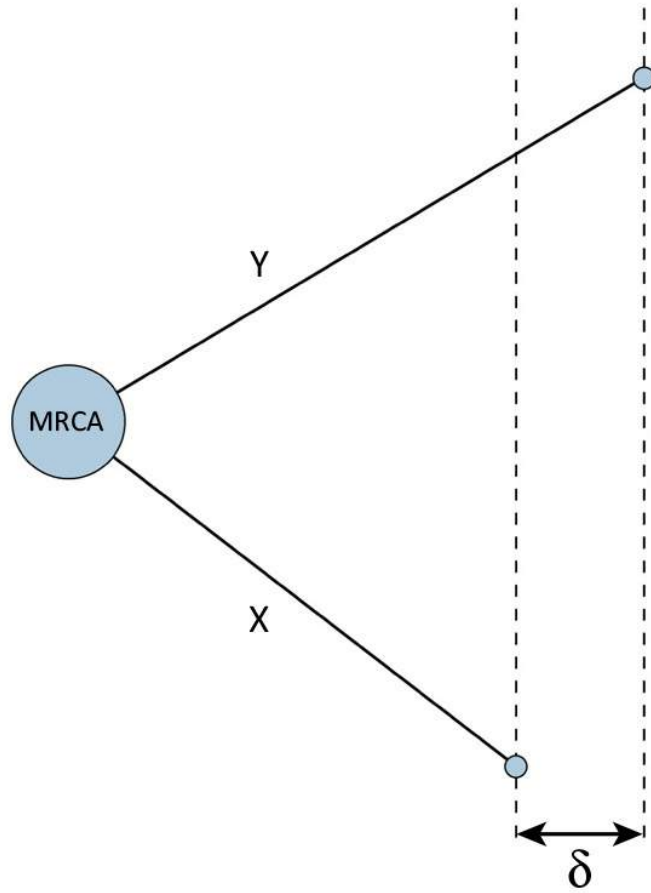
Estimating rates



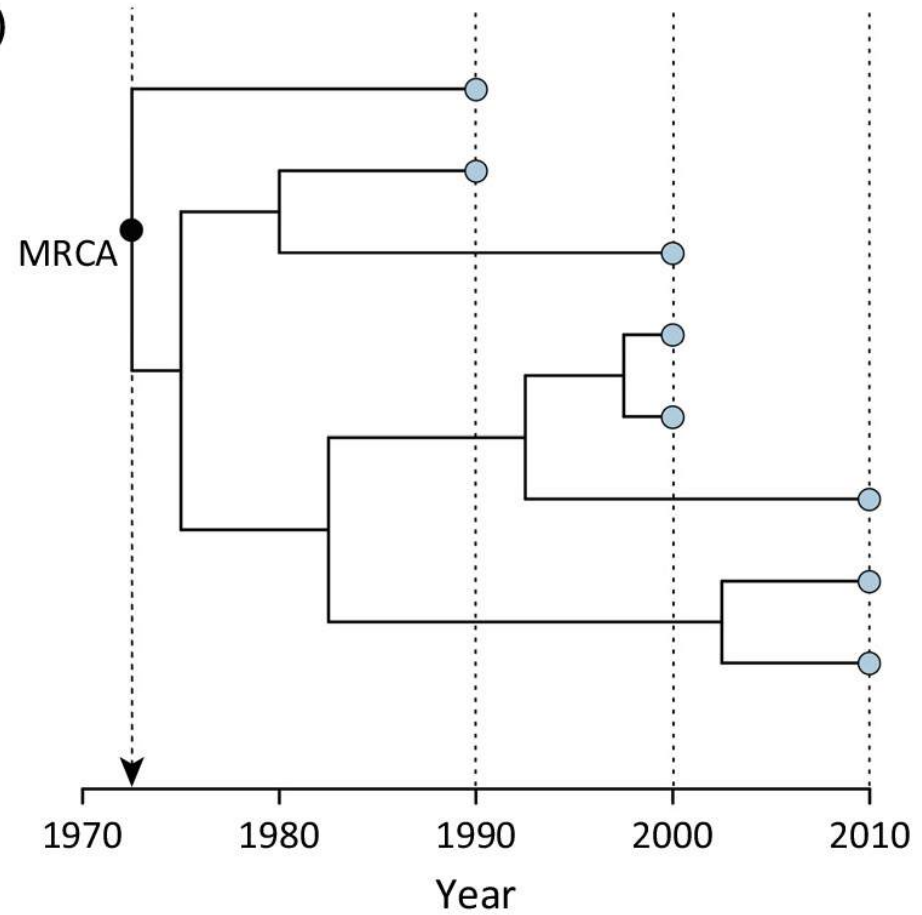
$$\frac{|dB - dA|}{|tA - tB|} = \text{Substitution rate}$$

Estimating rates

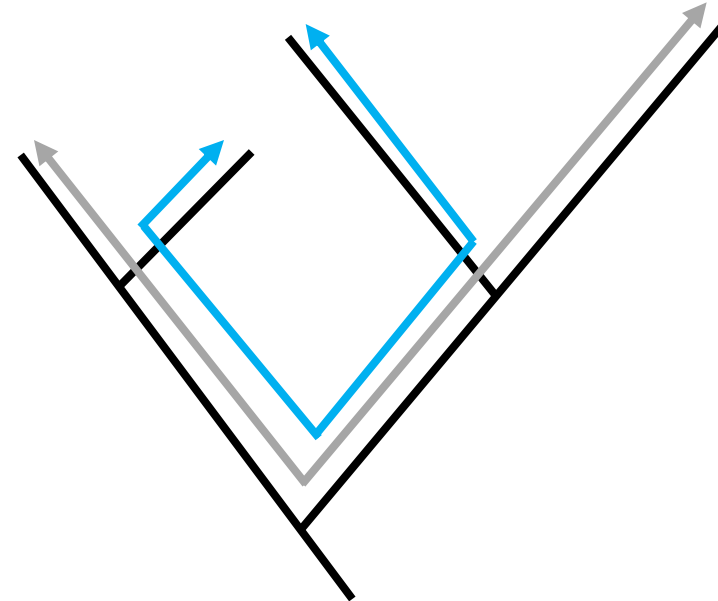
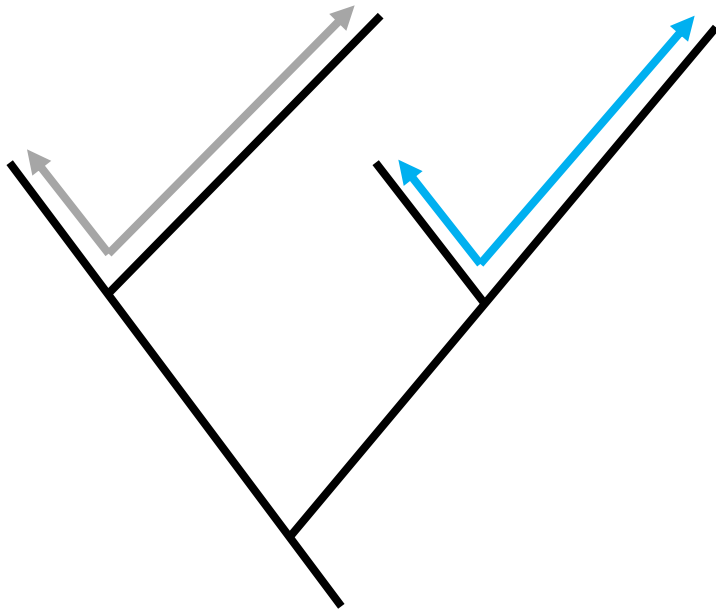
(A)



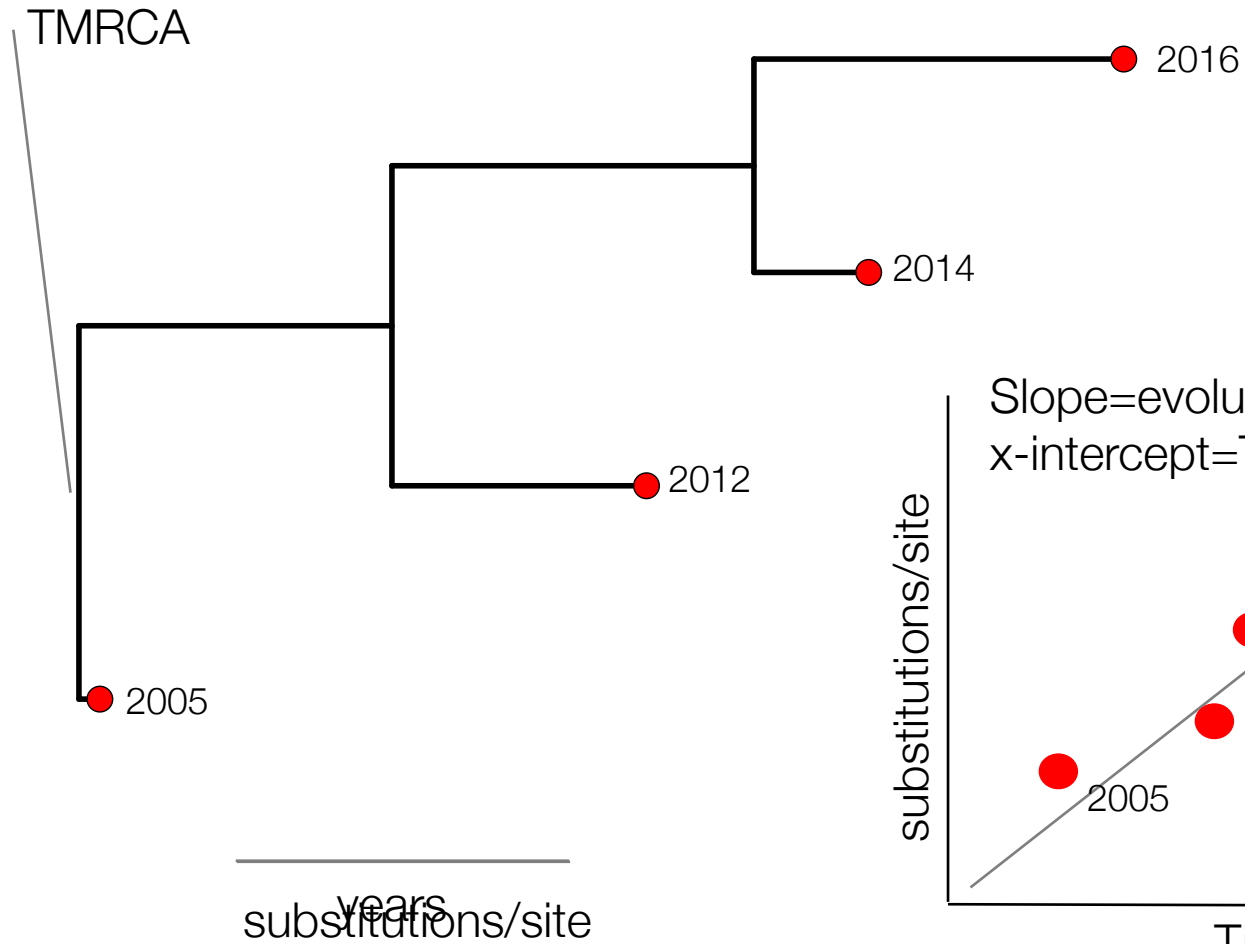
(B)



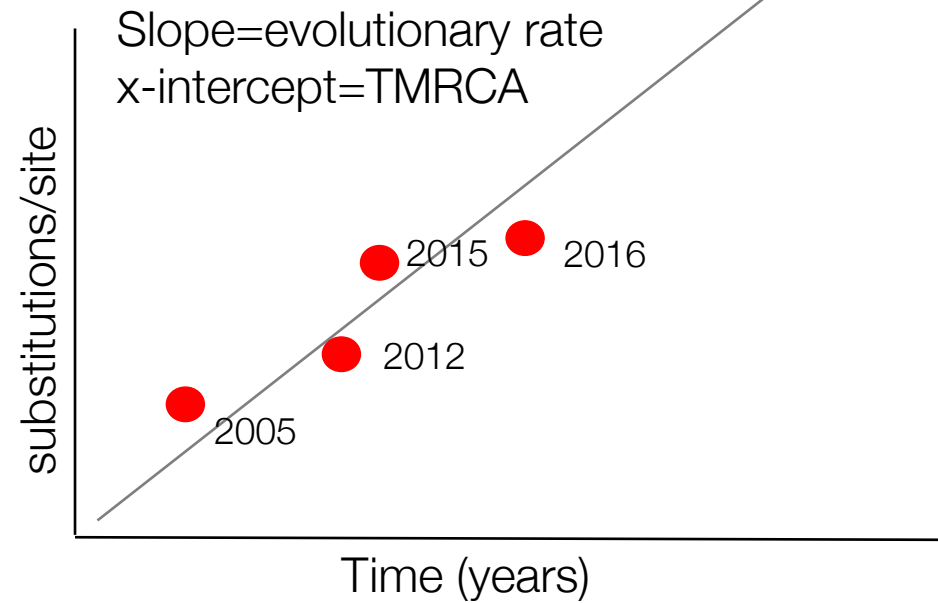
Estimating rates: phylo-temporal clustering

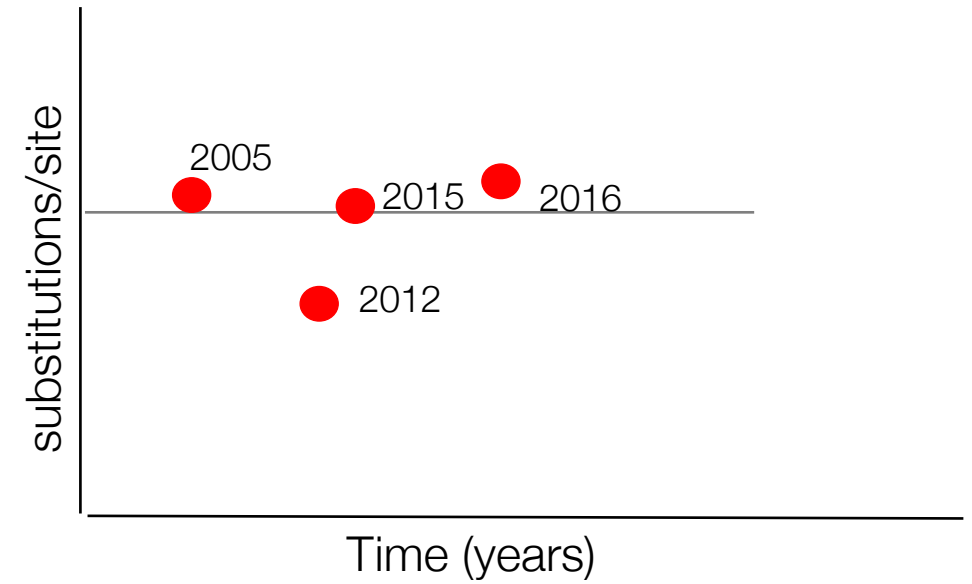
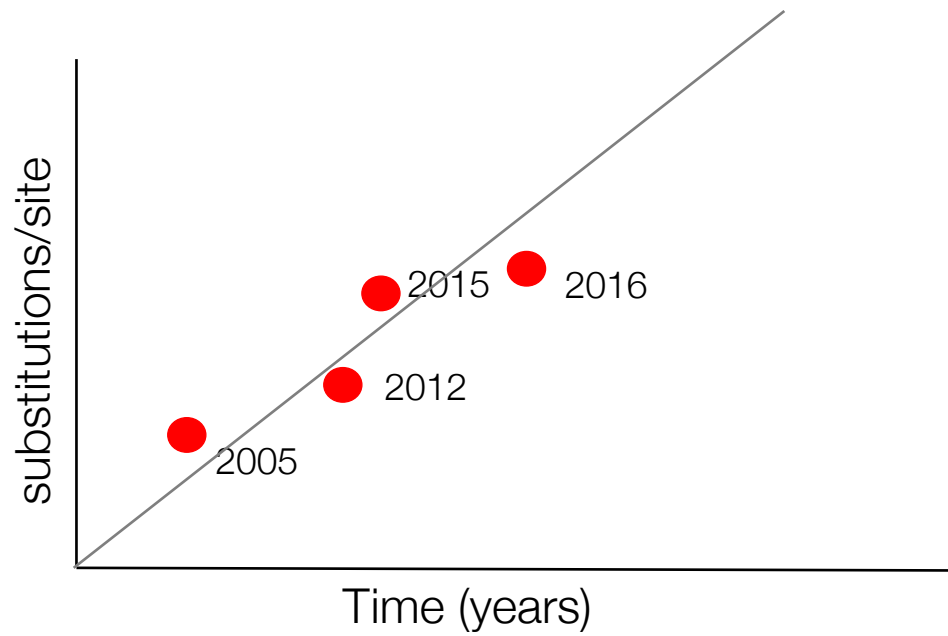


Assessing temporal structure

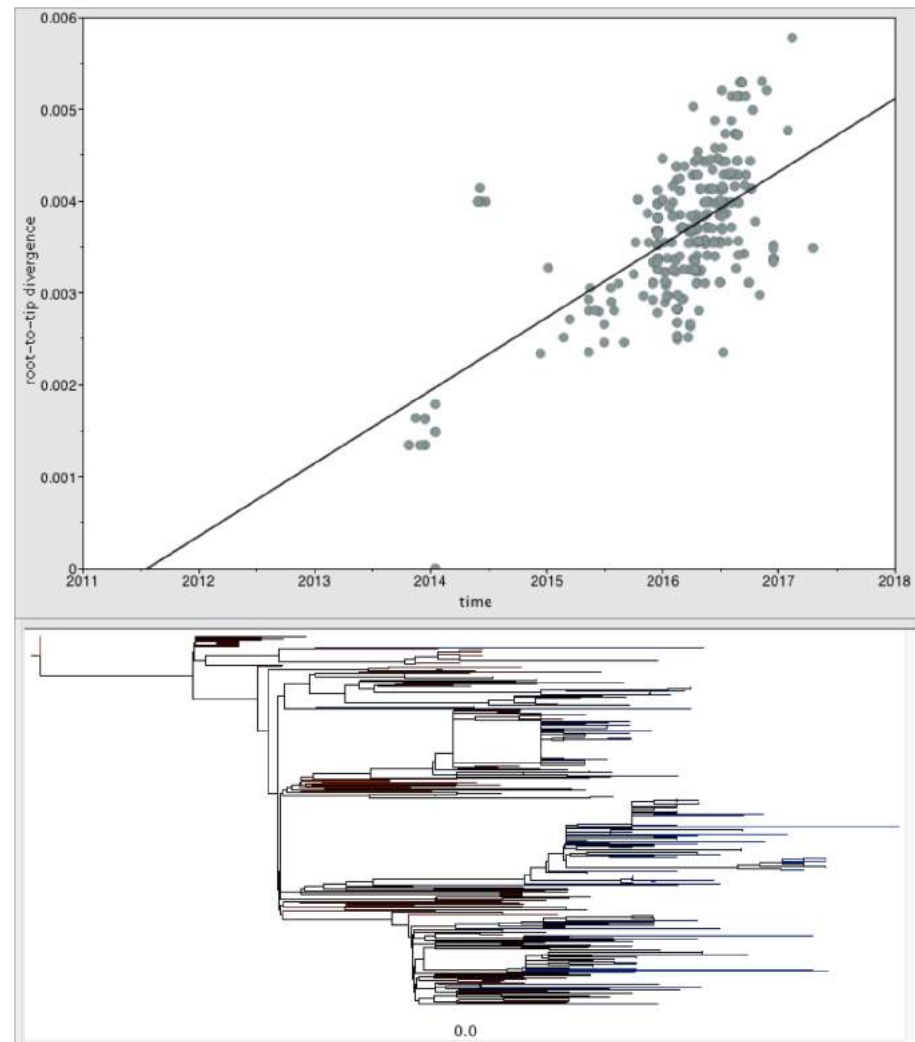
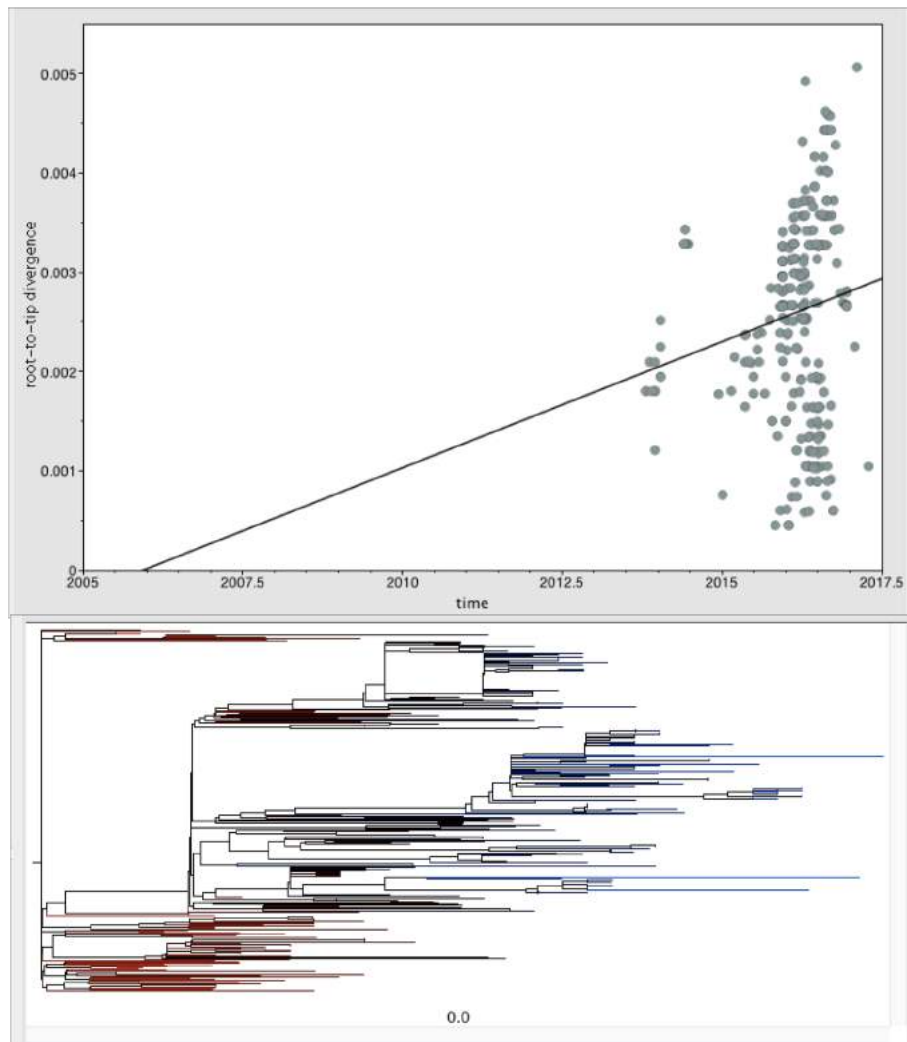


These estimates are
phylogenetically non-
independent!

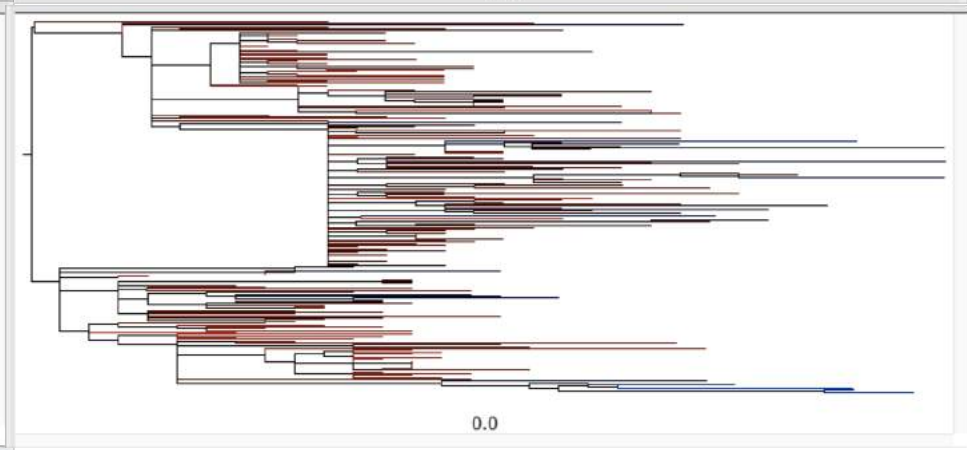
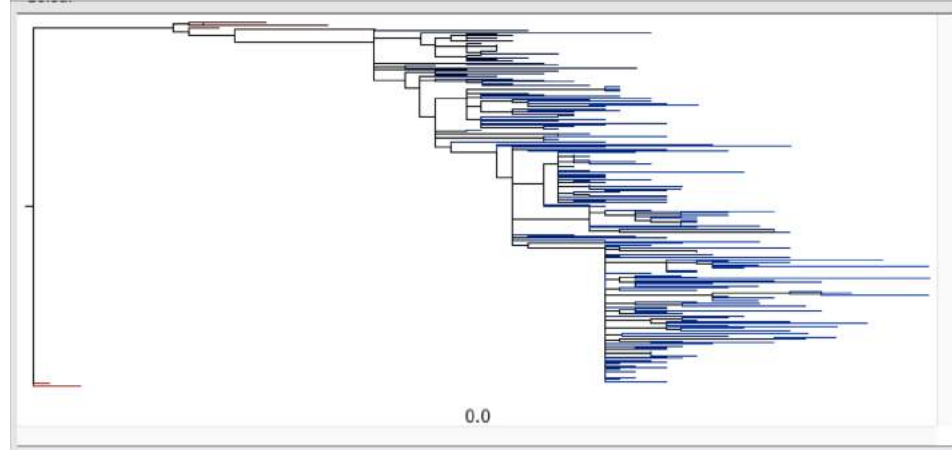
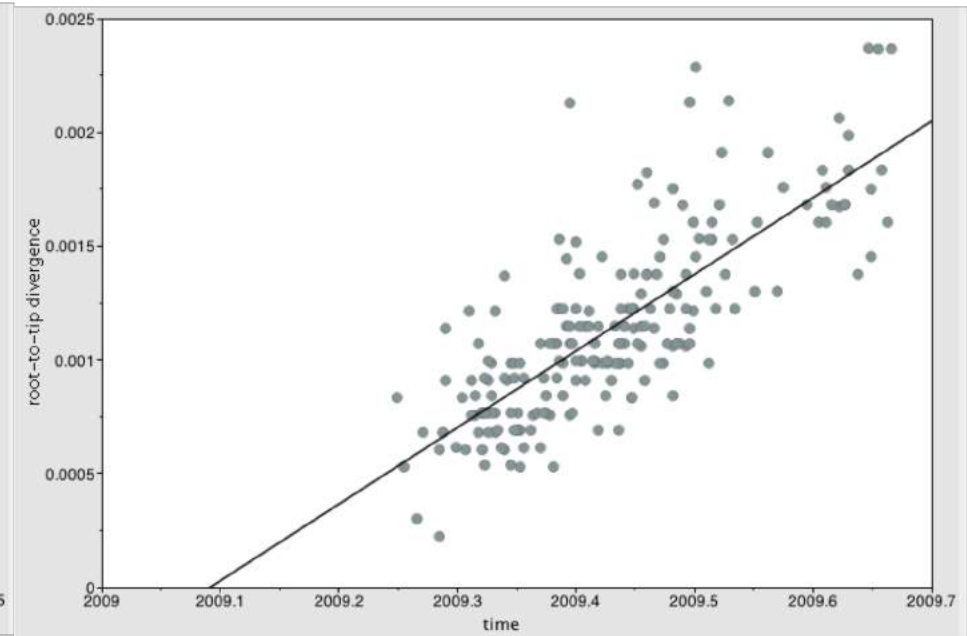
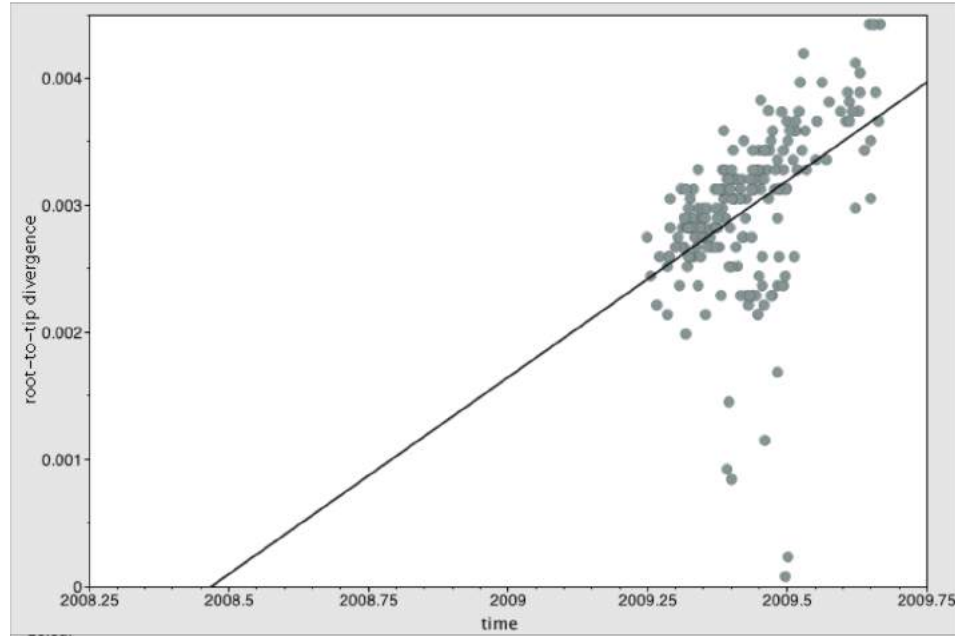




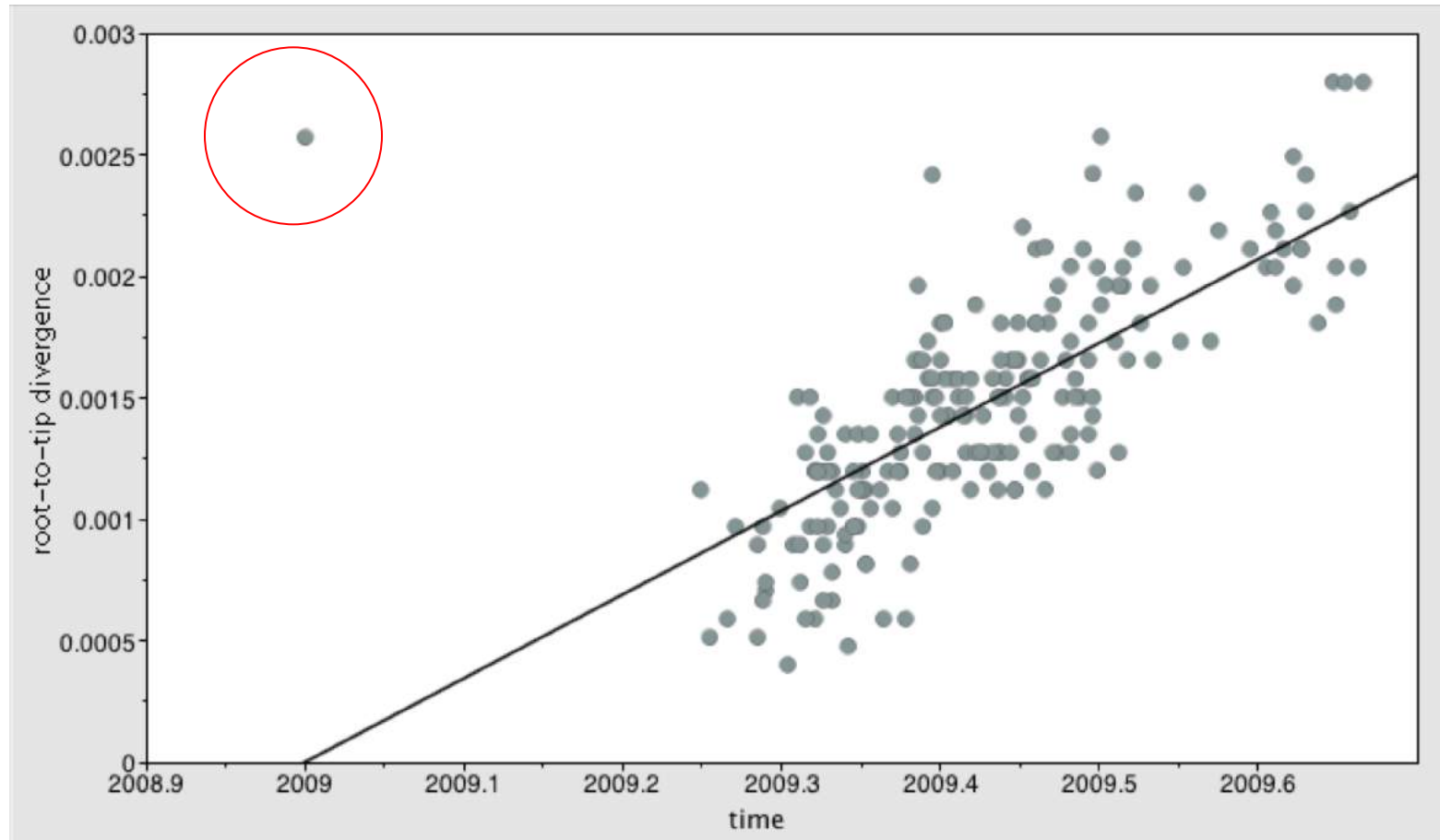
Slope should be positive
High R^2 (i.e. clocklike behavior)
Do not use p-value!



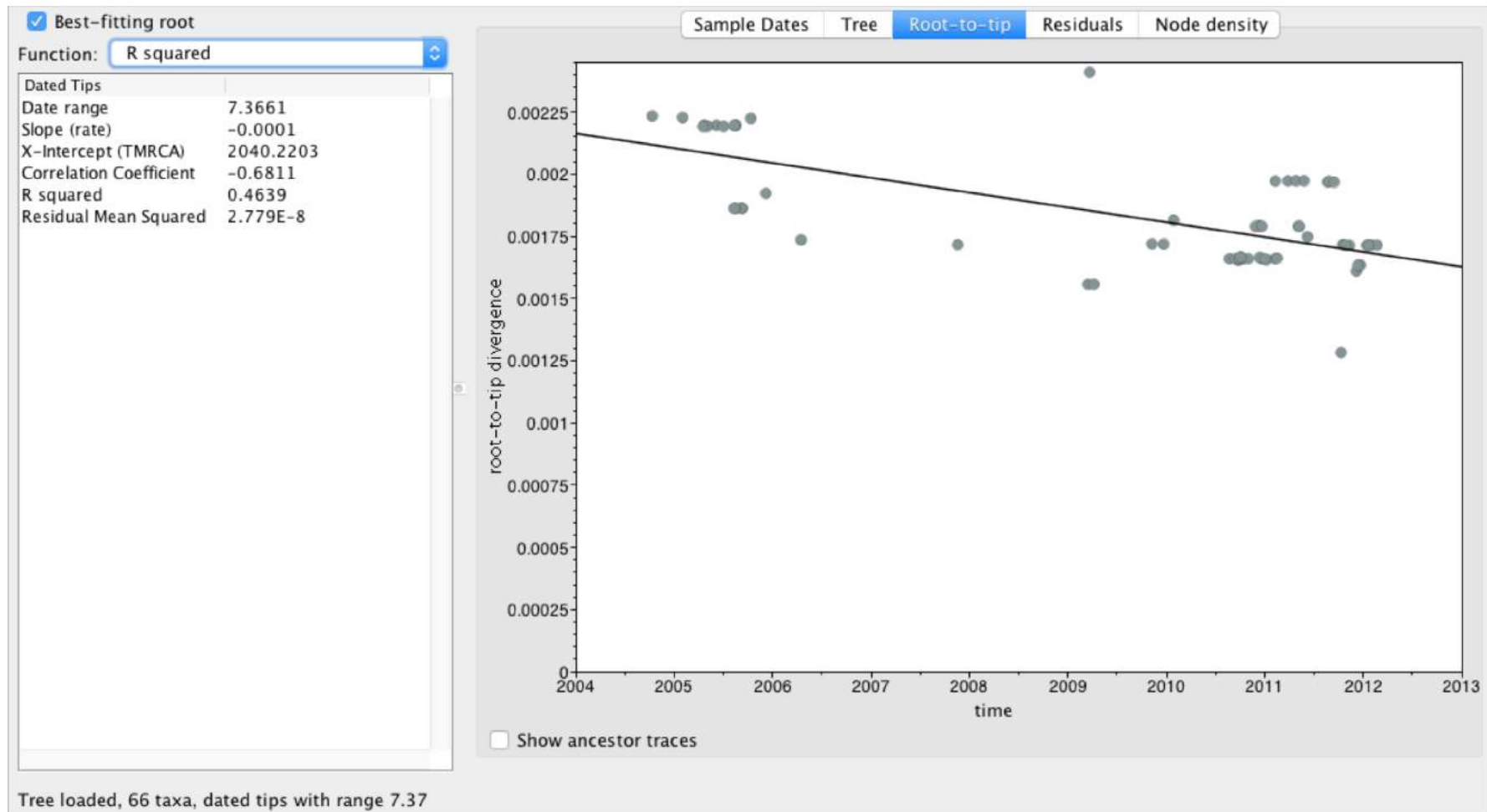
Zika virus in the Americas



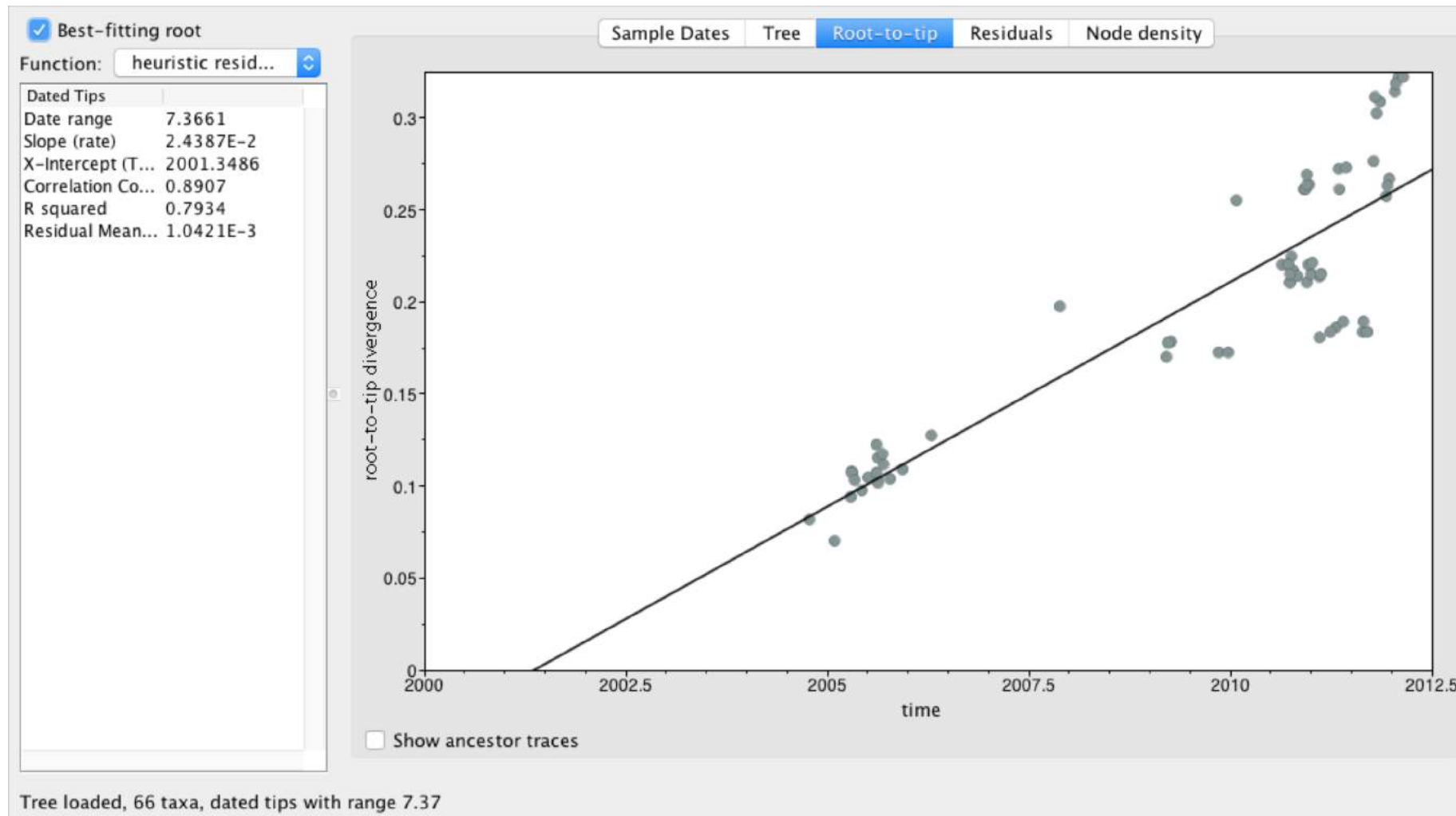
H1N1 flu



Outliers?

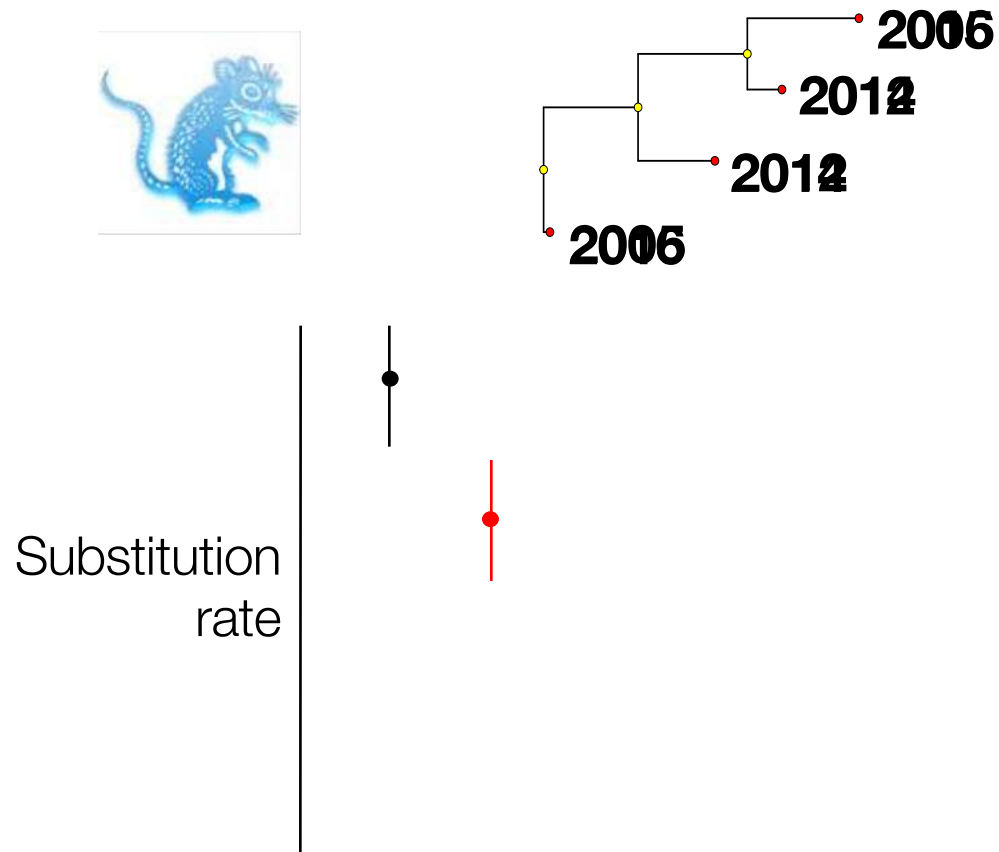


Acinetobacter baumannii

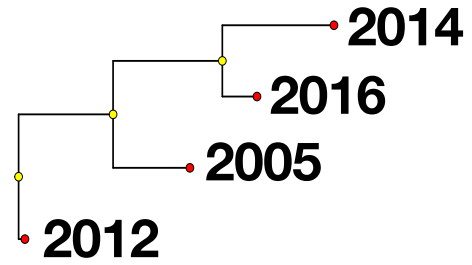


Acinetobacter baumannii

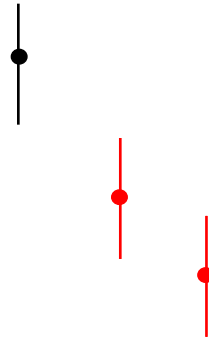
The date-randomisation test



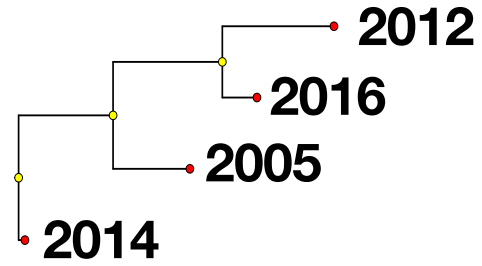
The date-randomisation test



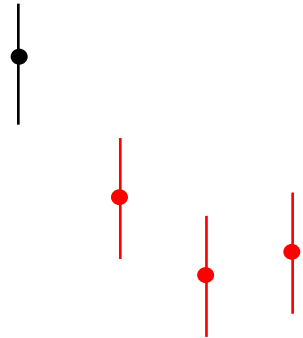
Substitution
rate



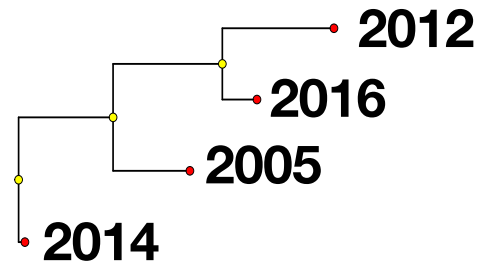
The date-randomisation test



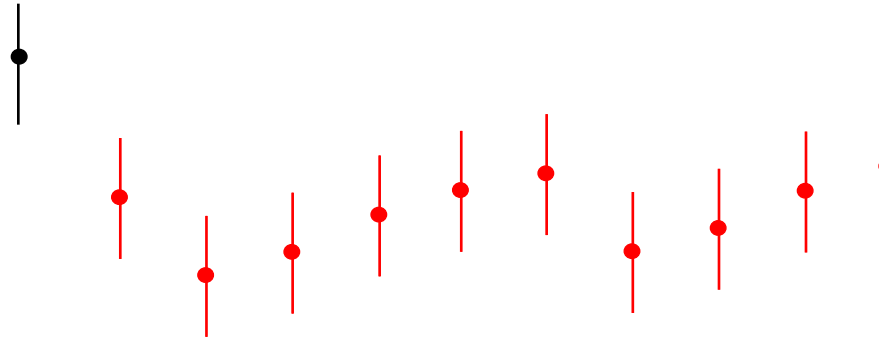
Substitution
rate



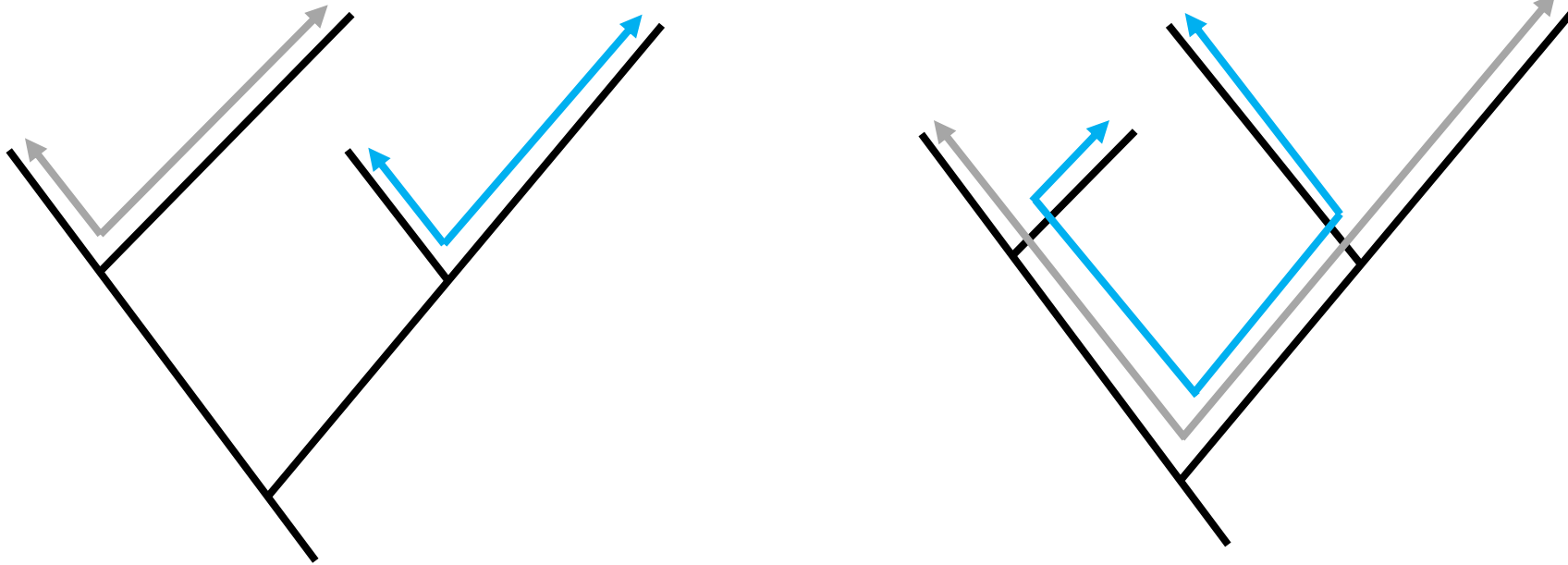
The date-randomisation test



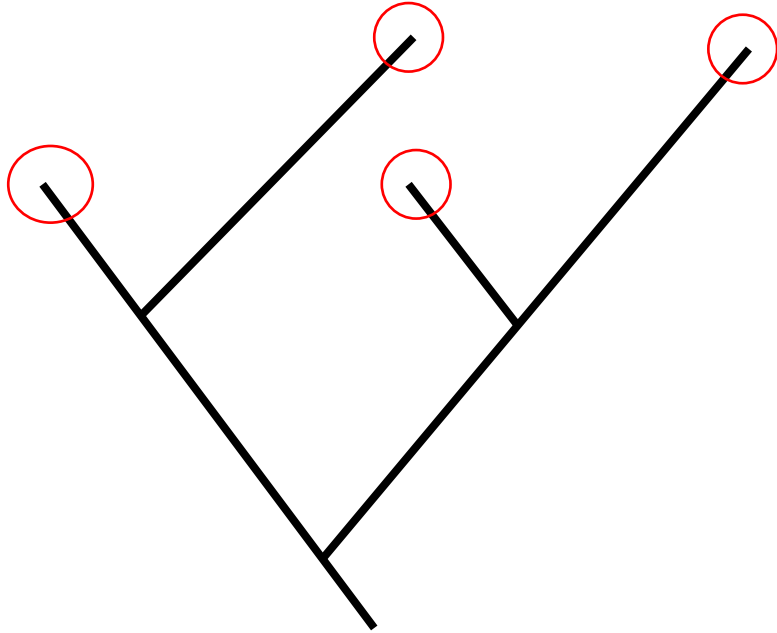
Substitution
rate



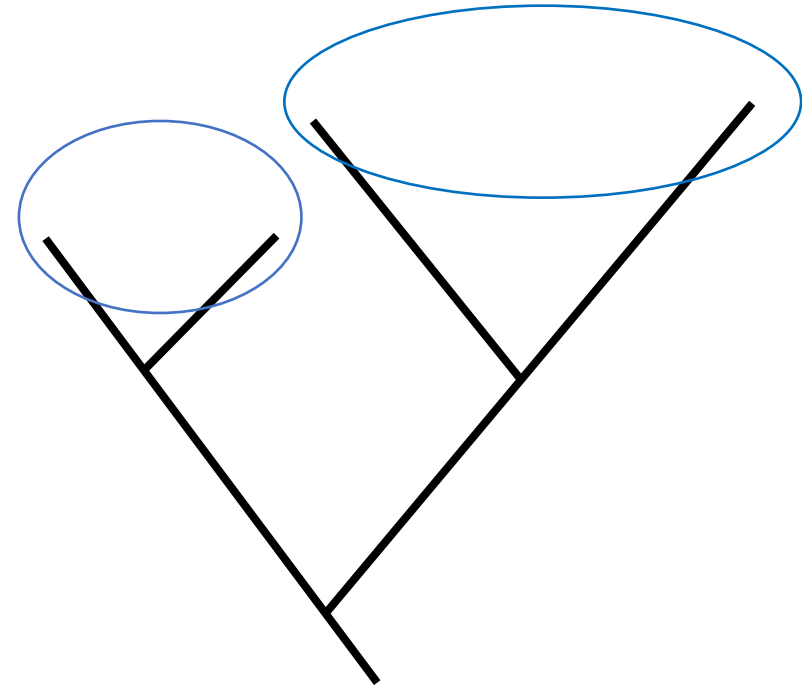
Phylogenetic and temporal clustering



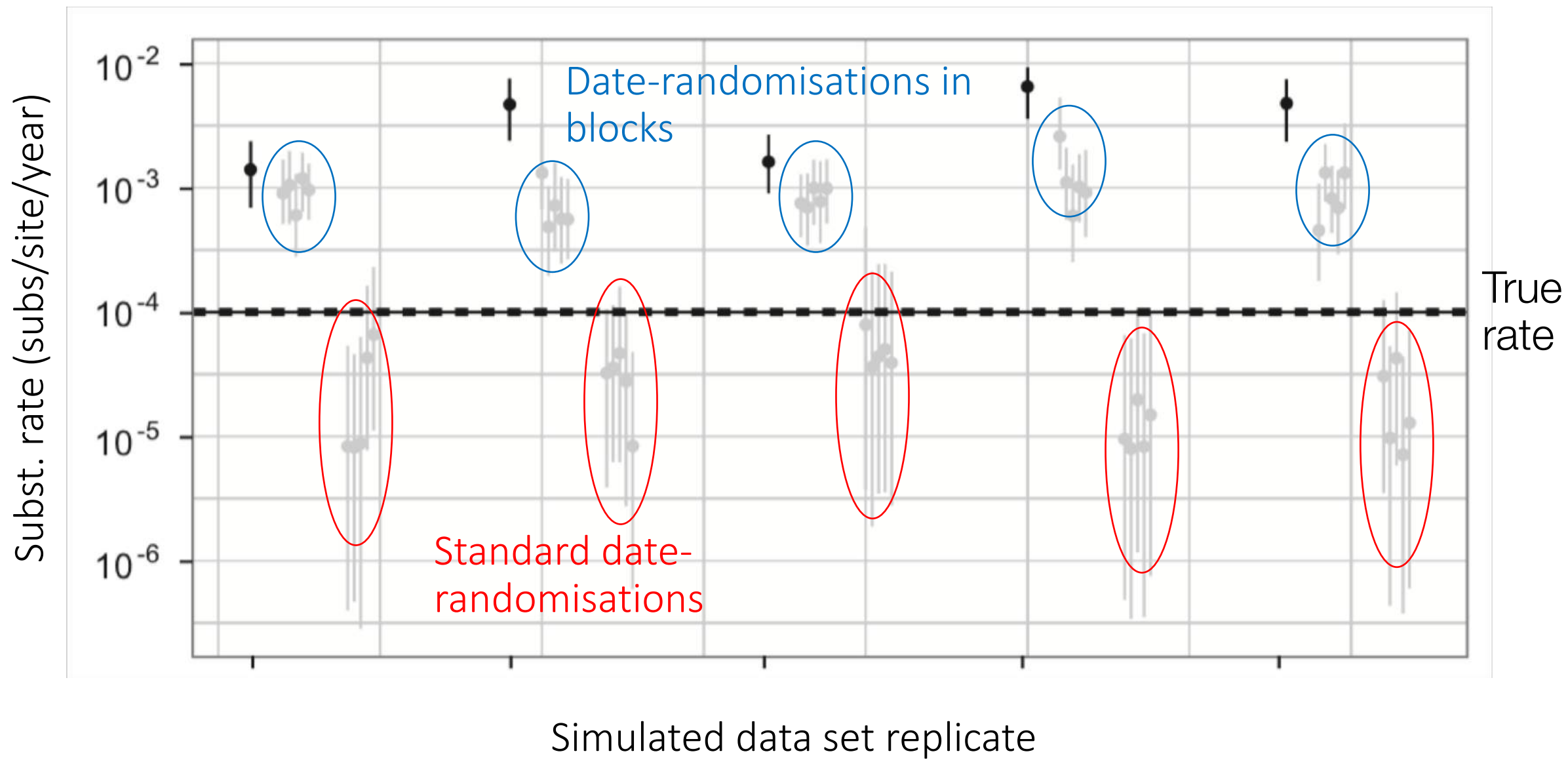
Date-randomisations in blocks



Murray *et al.* 2016 *Methods Ecol. Evol.*



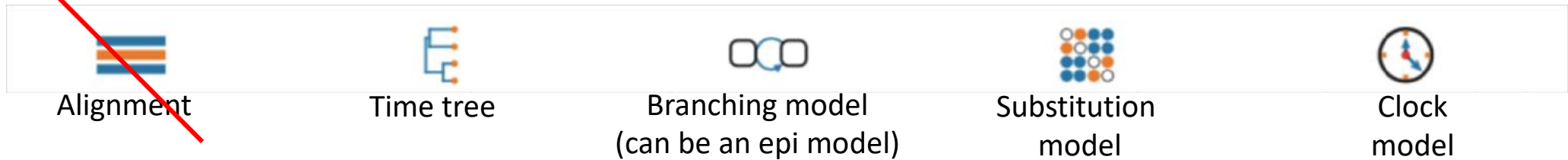
Duchêne *et al.* 2015 *Mol. Biol. Evol.*



Bayesian methods to assess temporal structure

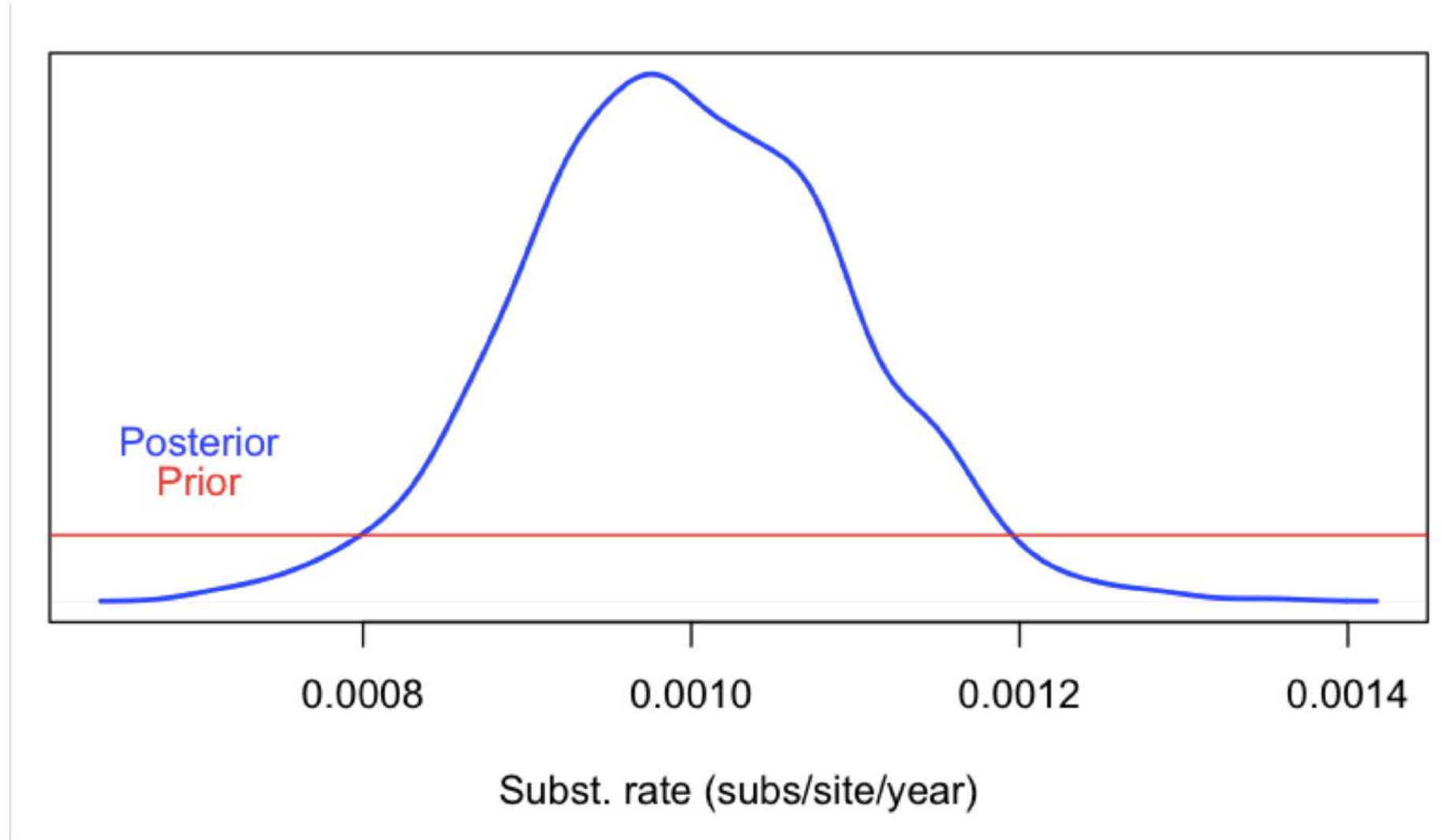
Comparing prior and posterior distributions

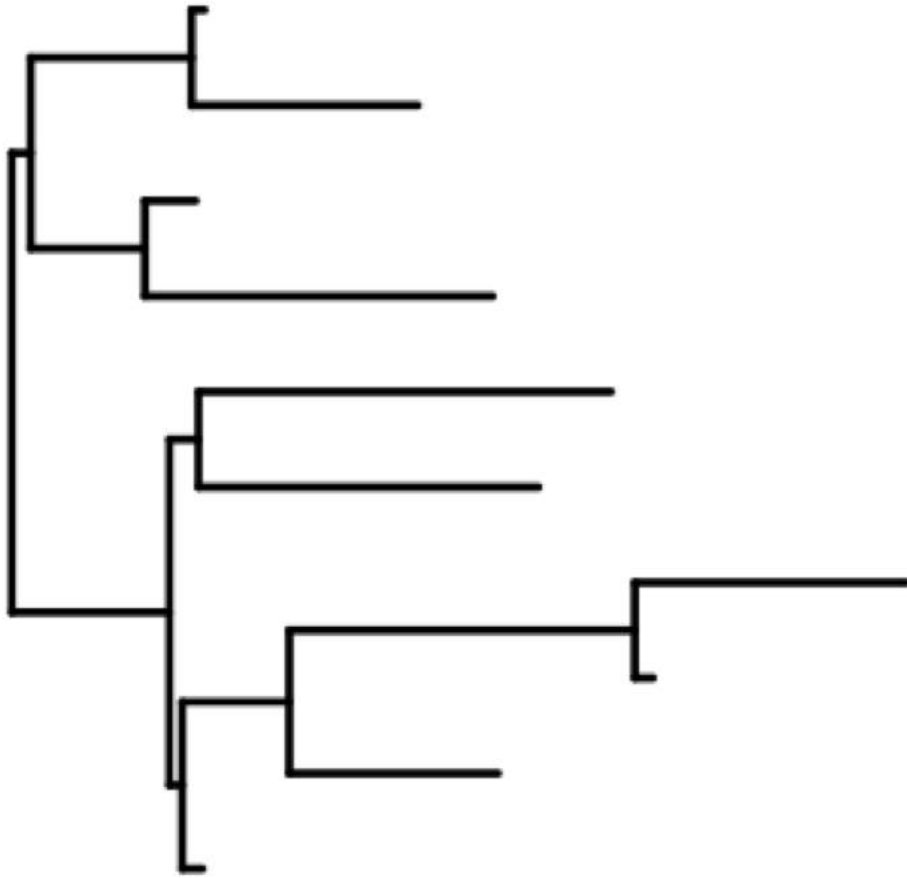
$$P(\text{E} \text{ } \text{ } \text{ } \text{ } | \text{ }) \propto P(\text{ } | \text{E} \text{ } \text{ } \text{ }) P(\text{E} | \text{ }) P(\text{ }) P(\text{ }) P(\text{ })$$



Sampling from the prior means that we do not calculate the phylogenetic likelihood.

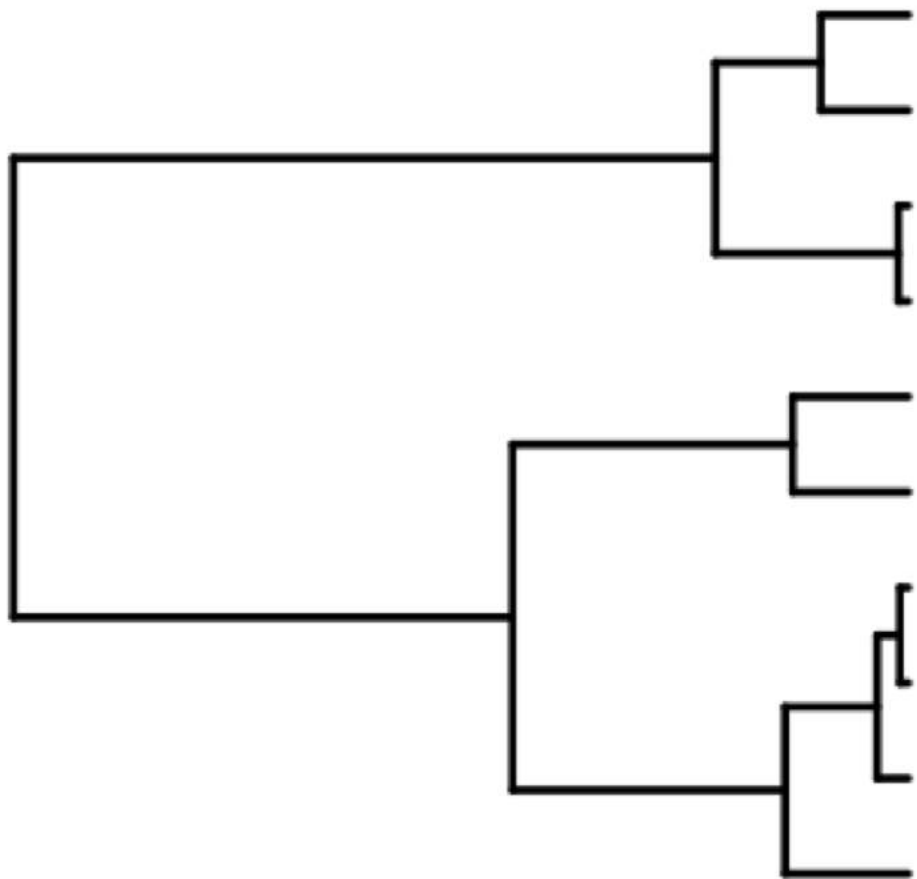
Comparing prior and posterior distributions





With no sequence data the root height and rate are driven by the prior.

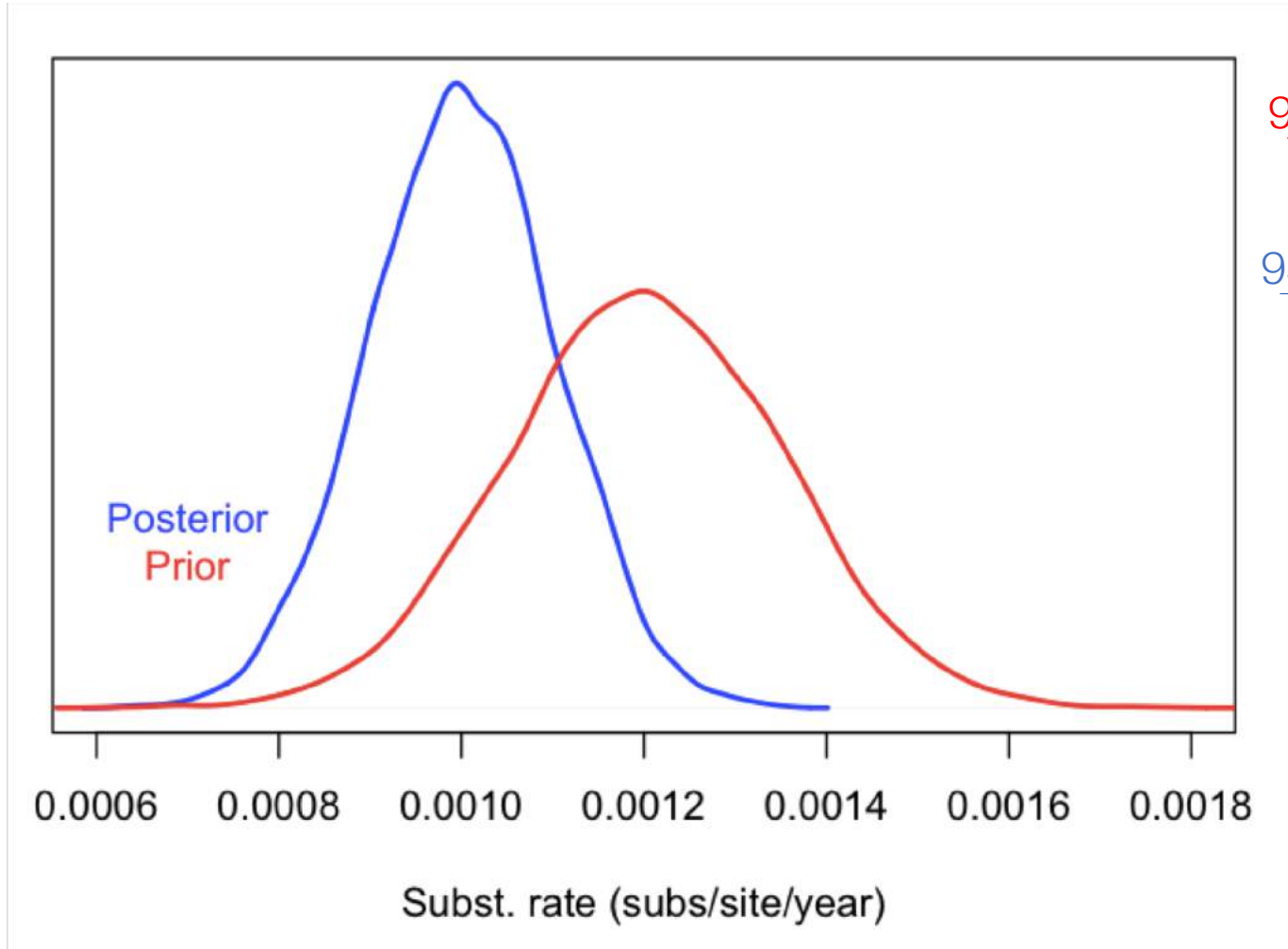
If the data have temporal structure the root height and rate should have a posterior that is much more informative than the prior.



If the sampling times are uninformative there is effectively no calibration information.

The root height should be driven by the prior.

Quantifying information content



$$\frac{\text{95\% quantile width}}{\text{Mean value}} = CV_{\text{prior}}$$

$$\frac{\text{95\% quantile width}}{\text{Mean value}} = CV_{\text{posterior}}$$

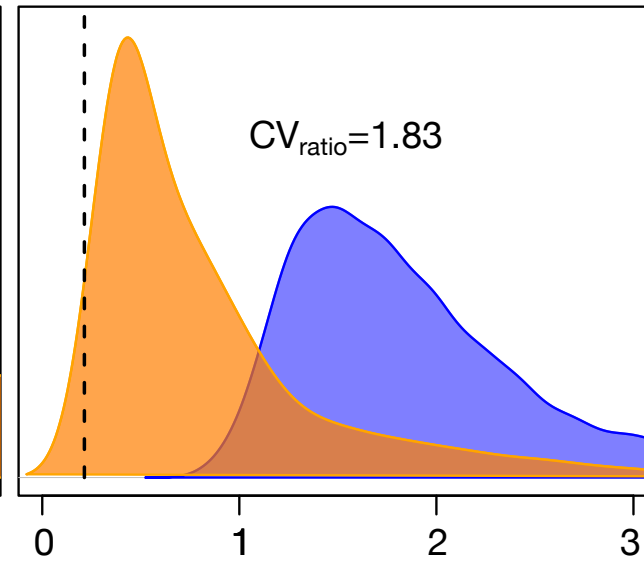
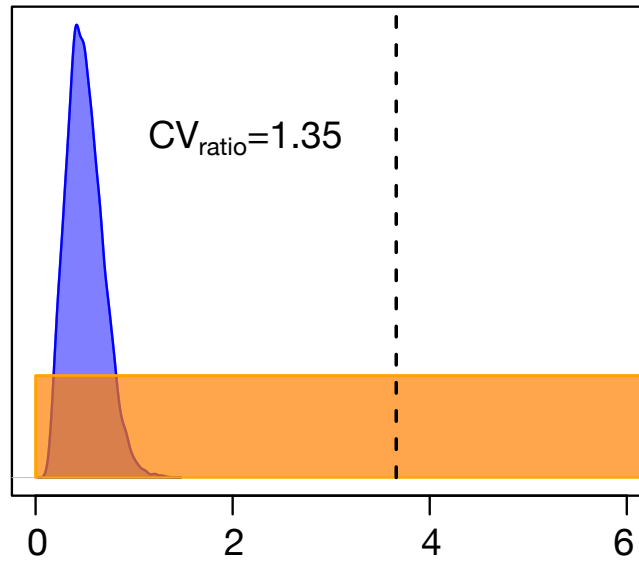
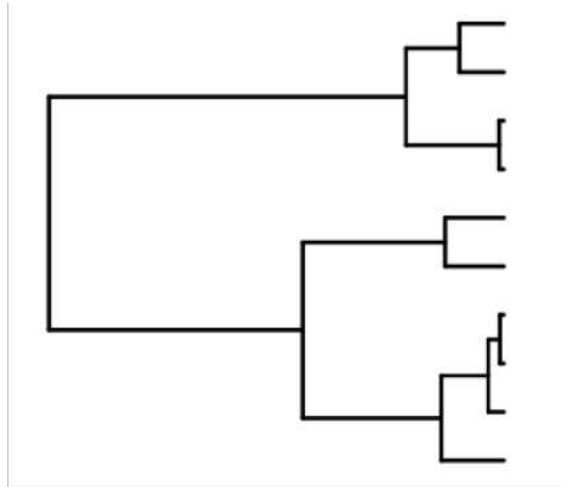
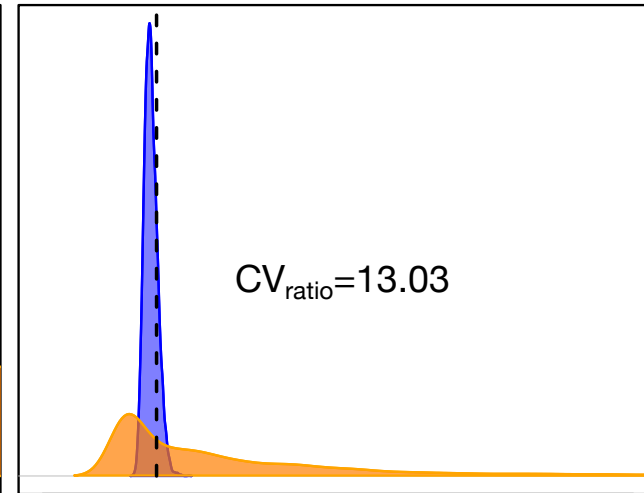
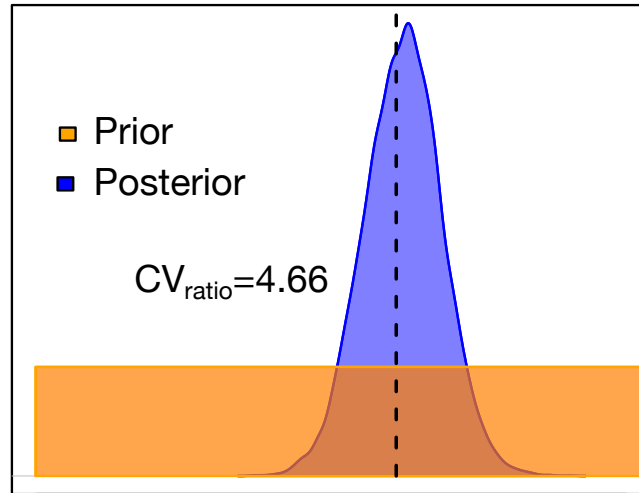
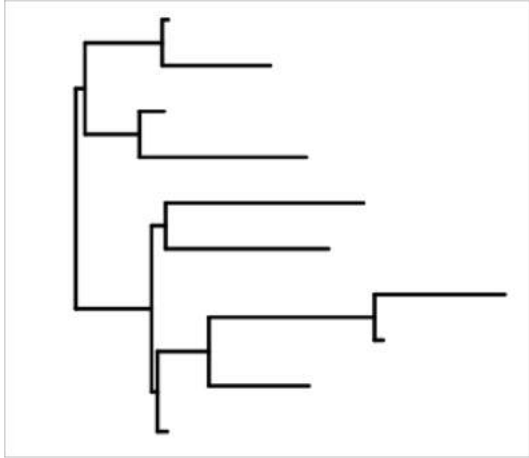
$$\frac{CV_{\text{prior}}}{CV_{\text{posterior}}} = CV_{\text{ratio}}$$

$CV_{\text{ratio}} > 1$ Posterior is more informative than the prior.

$CV_{\text{ratio}} = 1$ Posterior and prior are equally informative. i.e. estimates are driven by the prior.

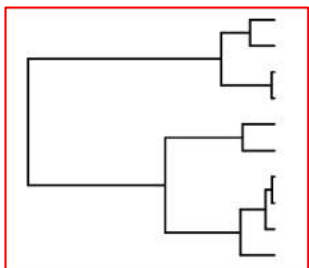
Evol. rate

Root height

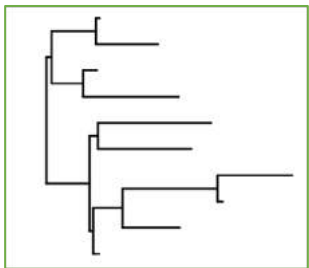


10^{-3} subs/site/year

Years

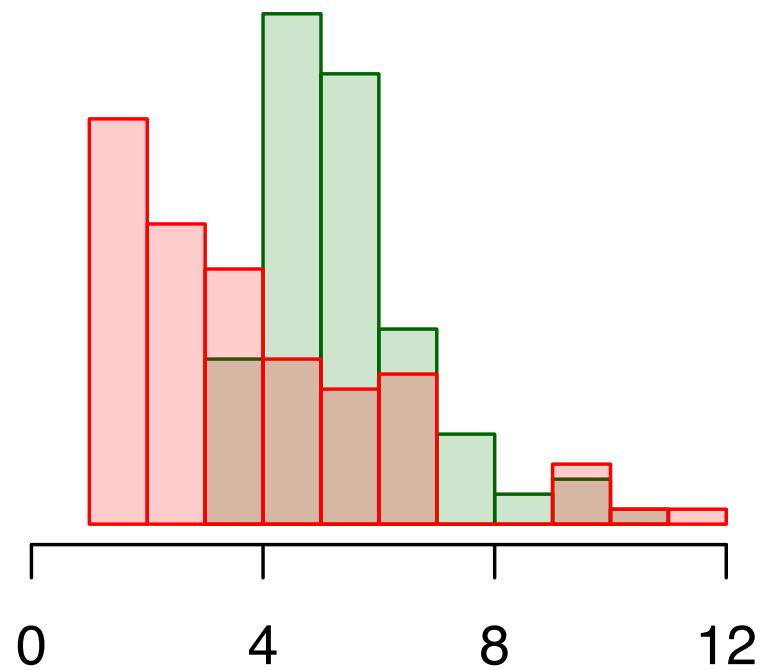


No temporal
structure

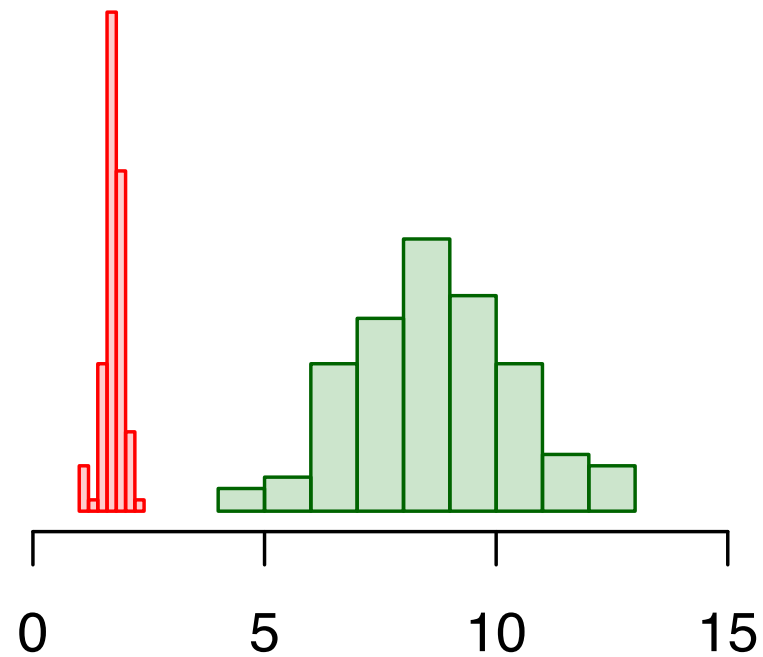


Temporal
structure

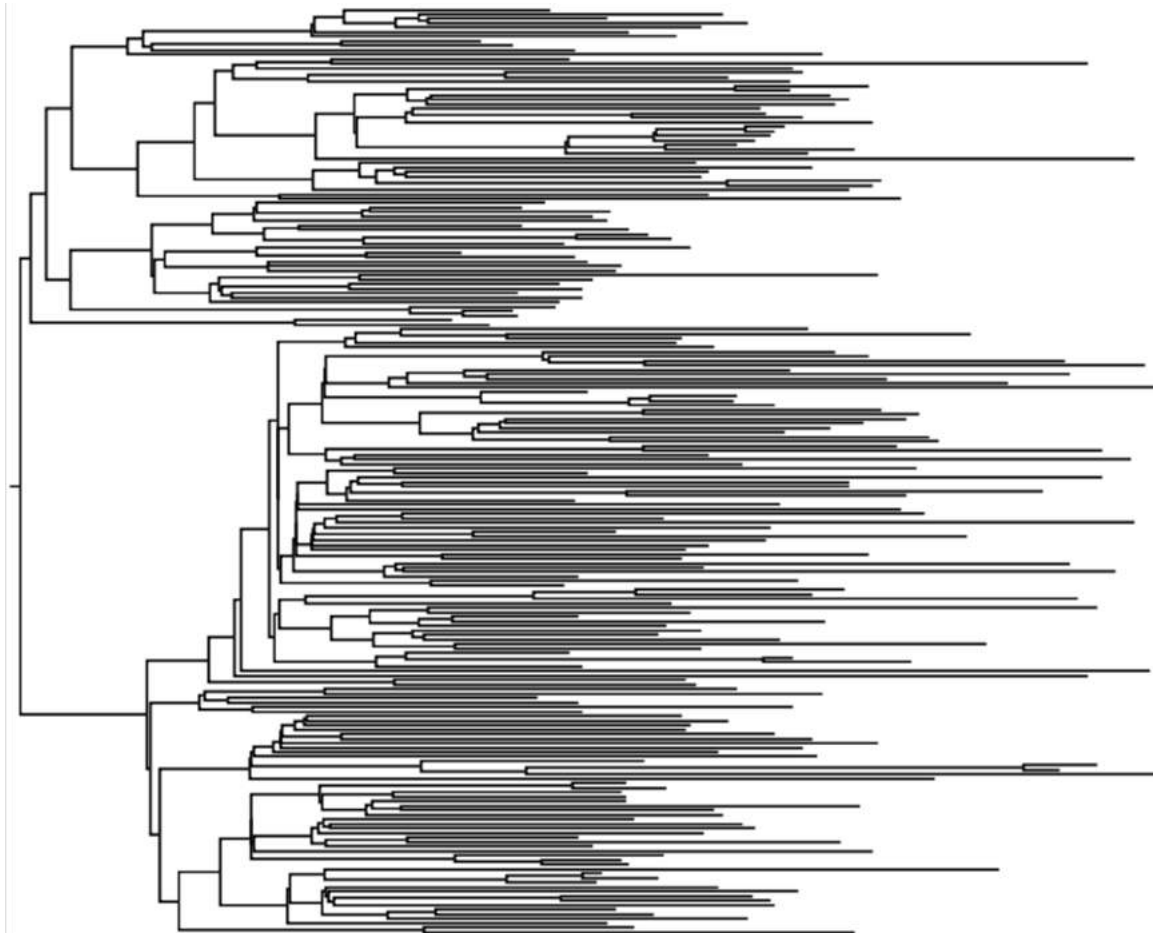
CV_{ratio} of Evol. rate



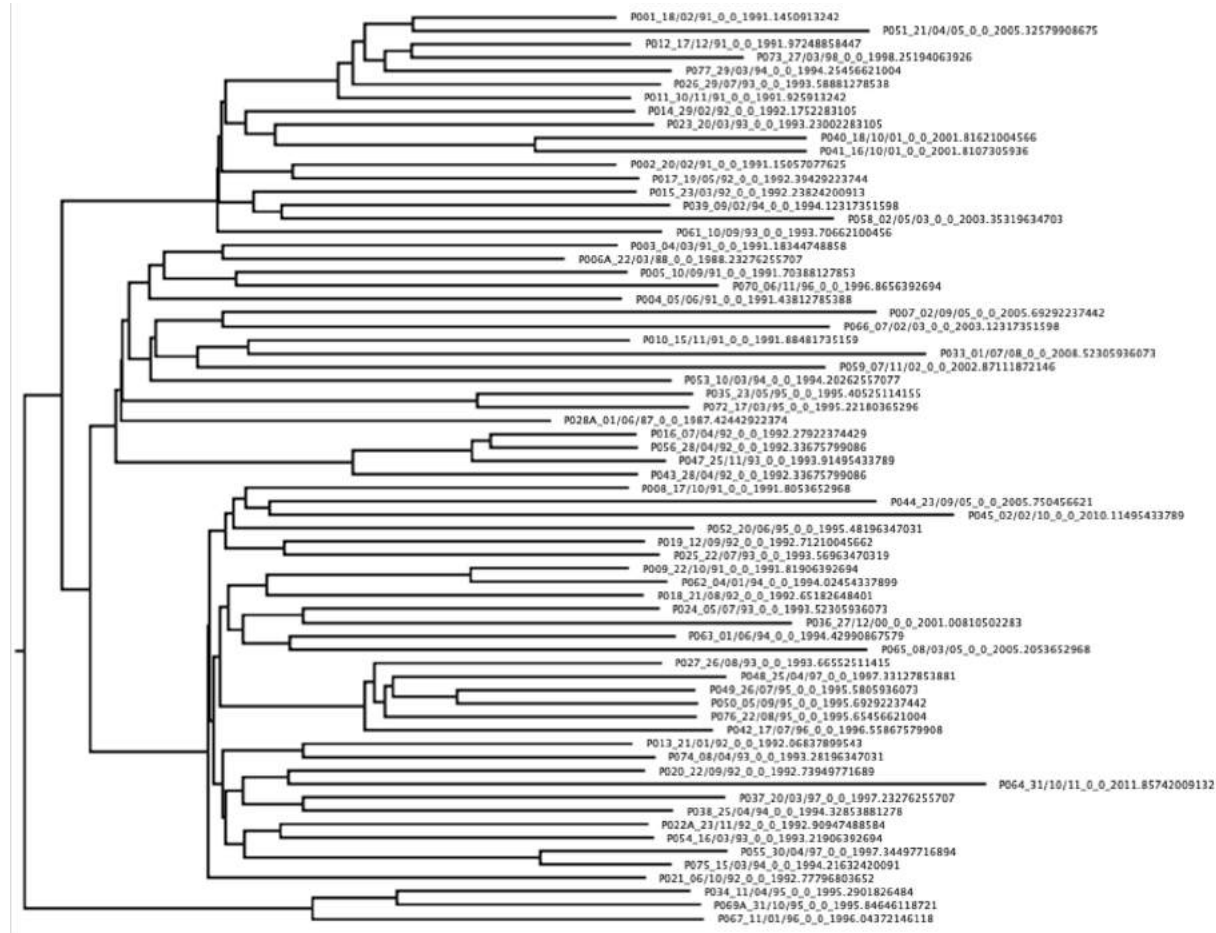
CV_{ratio} of root-height

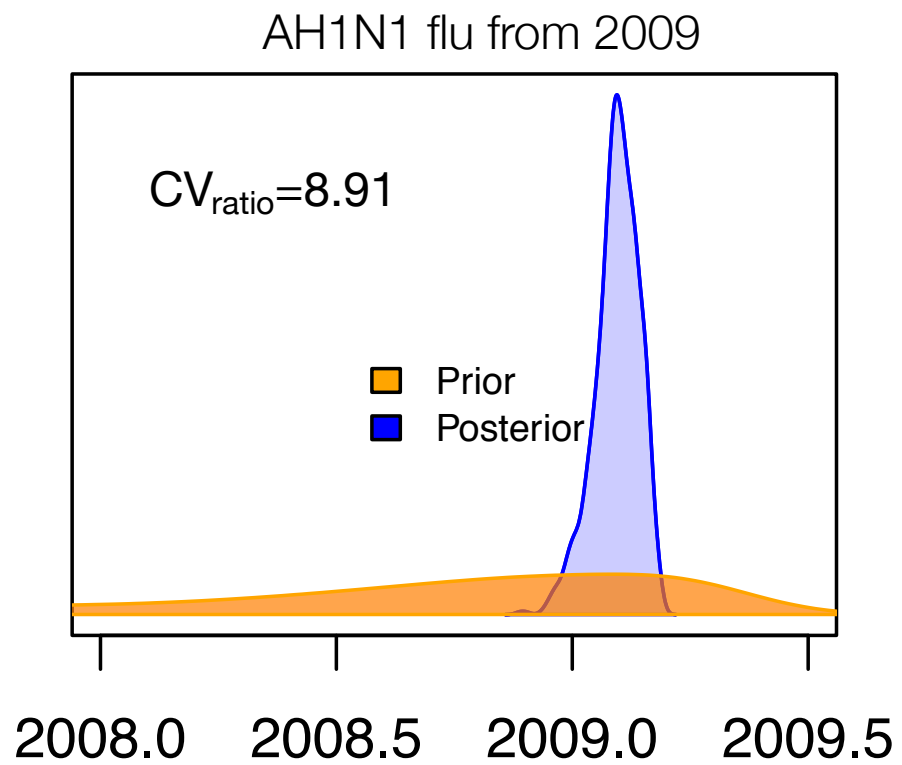


AH1N1 flu from 2009

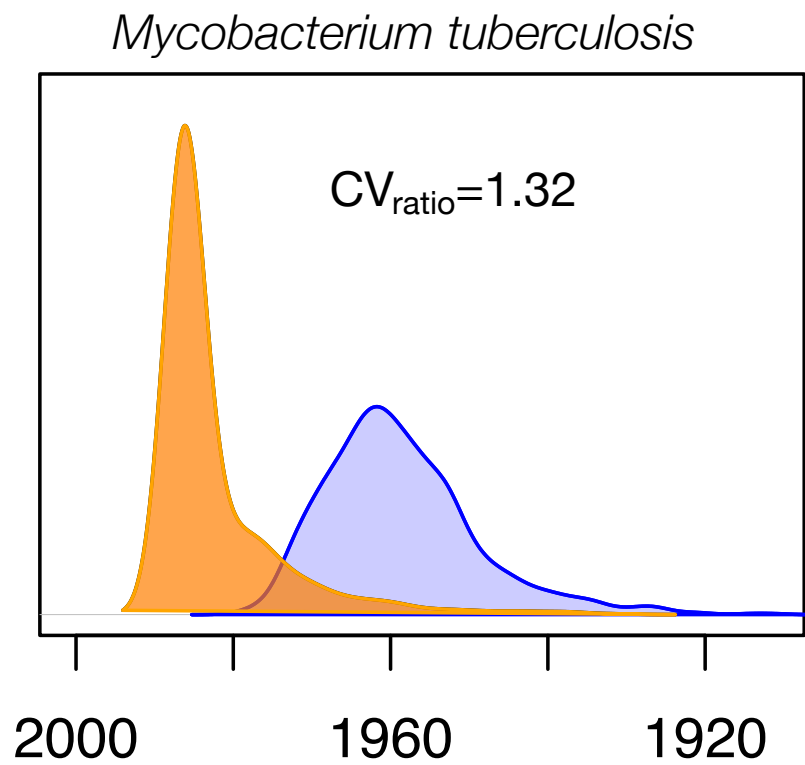


Mycobacterium tuberculosis





Age of root-node





Beast2

Bayesian evolutionary analysis by sampling trees

[ABOUT](#)[BOOK](#)[CITATION](#)[TUTORIALS](#)[FAQ](#)[BLOG](#)[DEVELOPERS ▾](#)

RECENT POSTS

What is new in v2.5.2

February 08, 2019

Assessing temporal structure
in data sets of rapidly-evolving
microbes and ancient DNA

January 28, 2019

Postdoc and two PhD posi-
tions in Auckland!

December 05, 2018

Workshop Announcement --

ASSESSING TEMPORAL STRUCTURE IN DATA SETS OF RAPIDLY-EVOLVING MICROBES AND ANCIENT DNA

Duchene and Duchene. Submitted.

Taming the BEAST



[news](#) [workshops](#) [tutorials](#) [contribute](#) 

Molecular dating using heterochronous data and substitution model averaging

Estimating the time of origin of the 2009 H1N1 pandemic in North America

by Sebastian Duchene , Tim Vaughan and Veronika Boskova

[Tutorial](#)

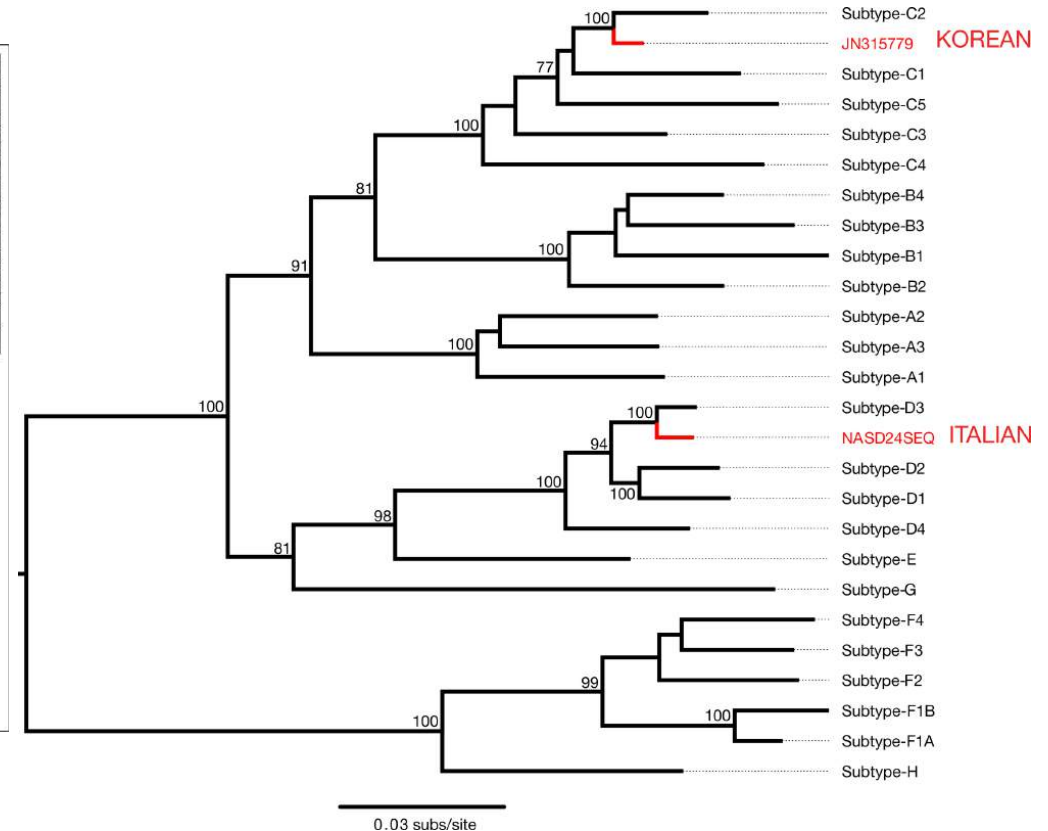
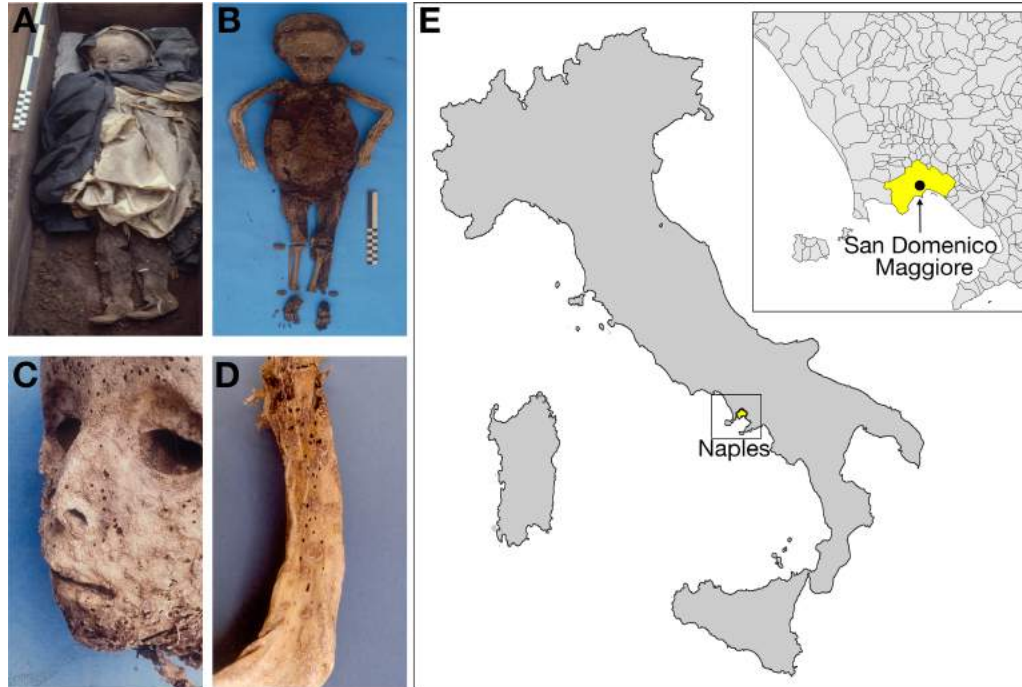
 [Github repository](#)

[Background](#)

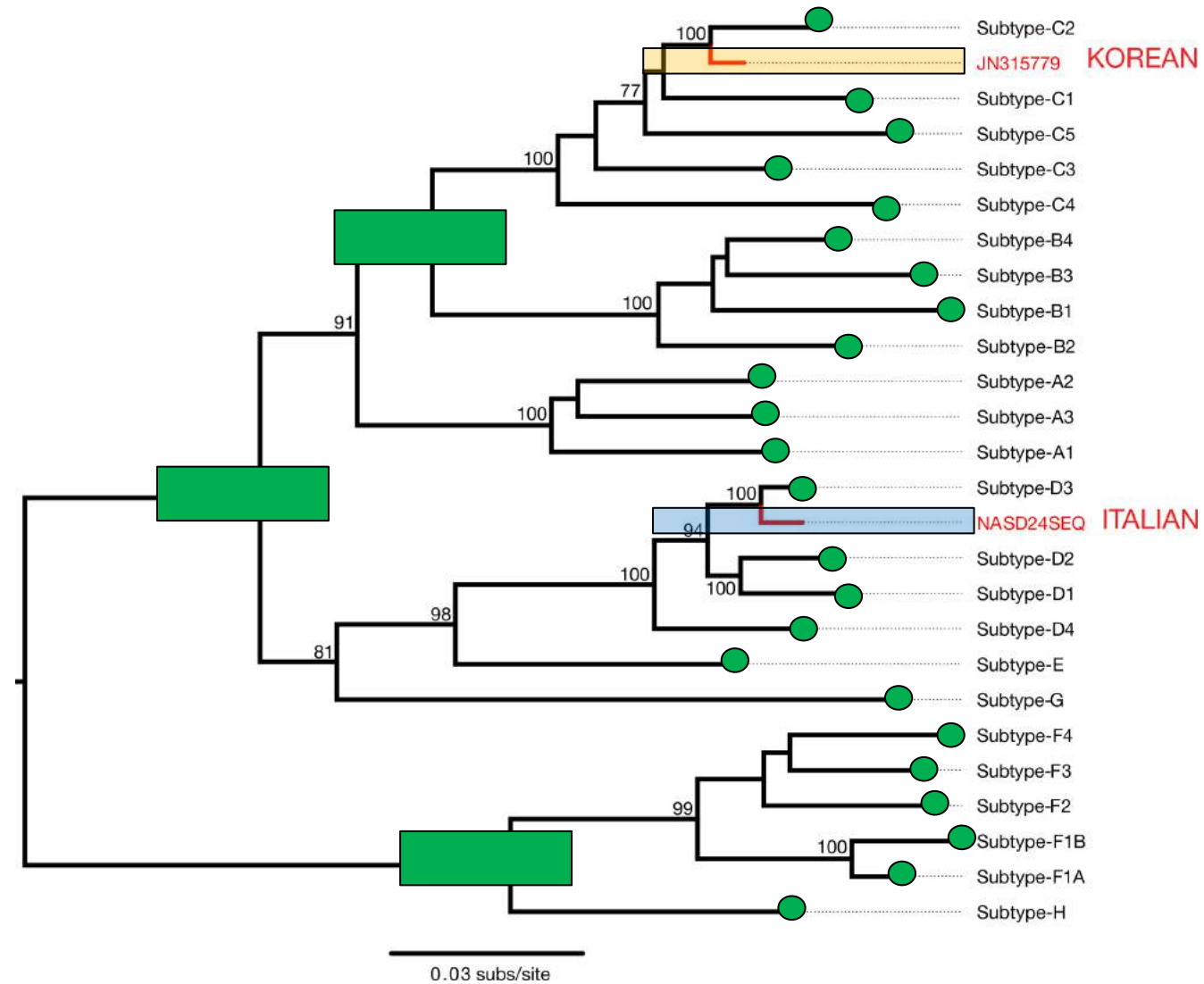
A case study in HBV (ancient DNA)

Comparing the prior and posterior -tips

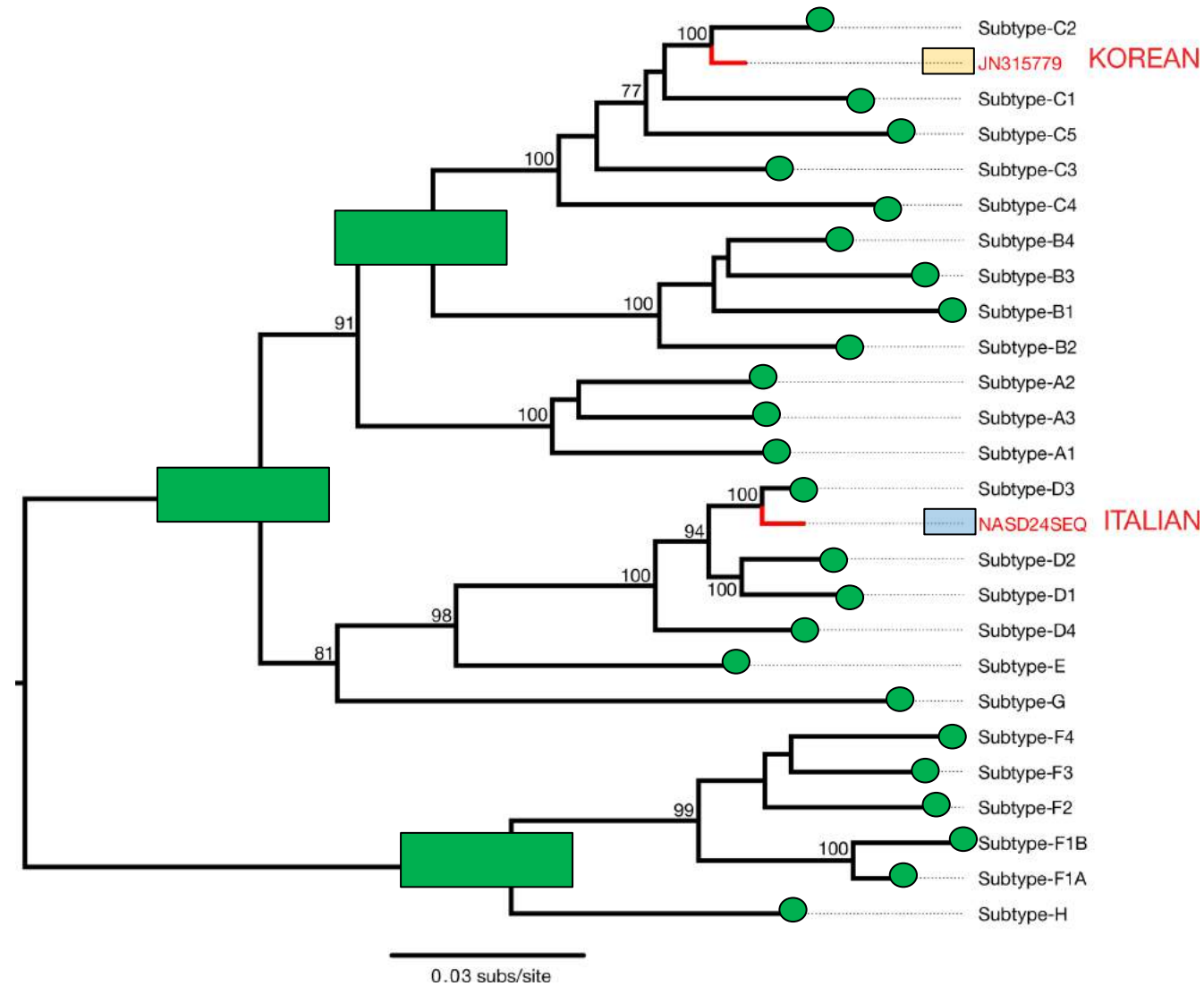
Hepatitis B Virus paradox



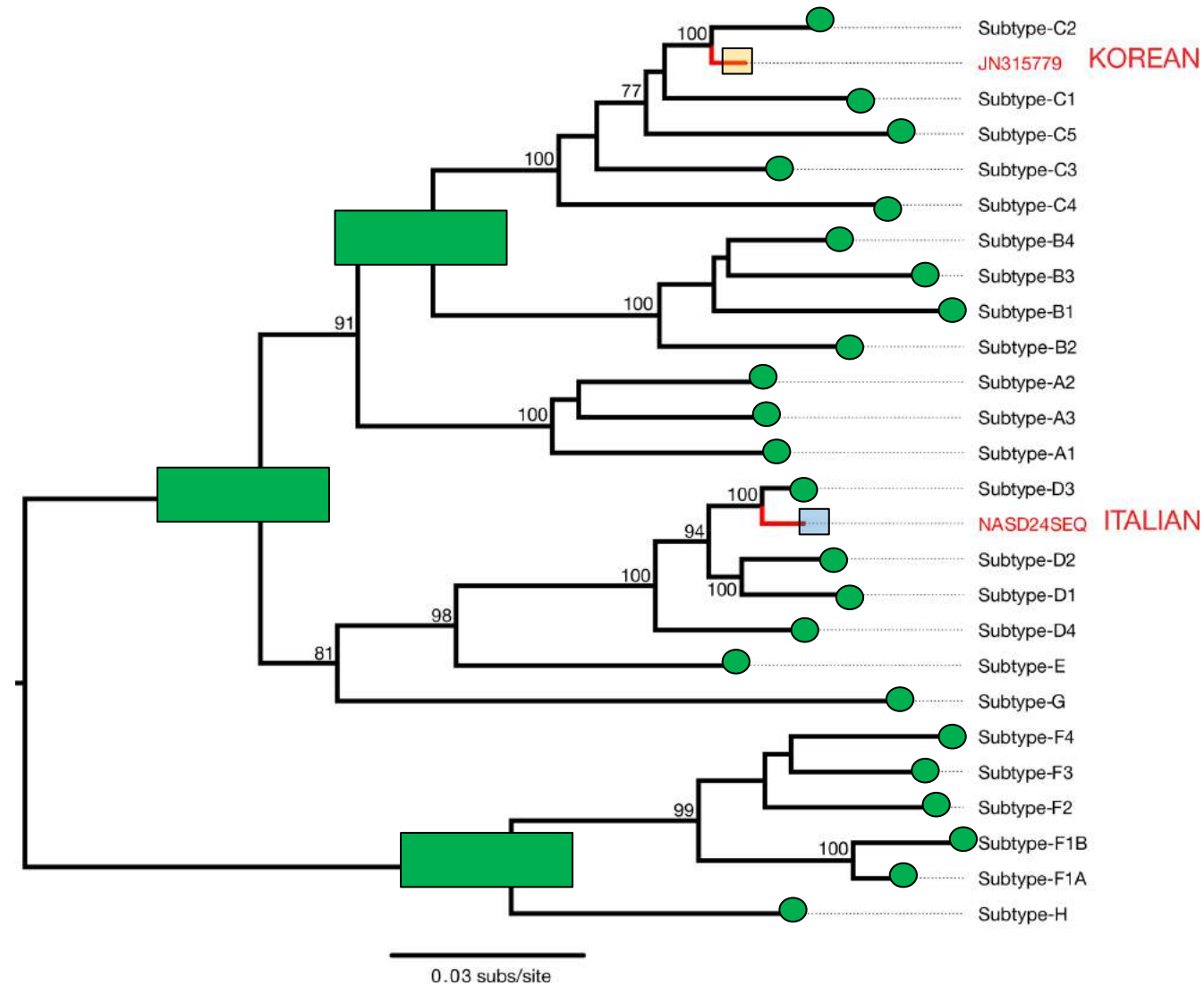
Comparing the prior and posterior -tips

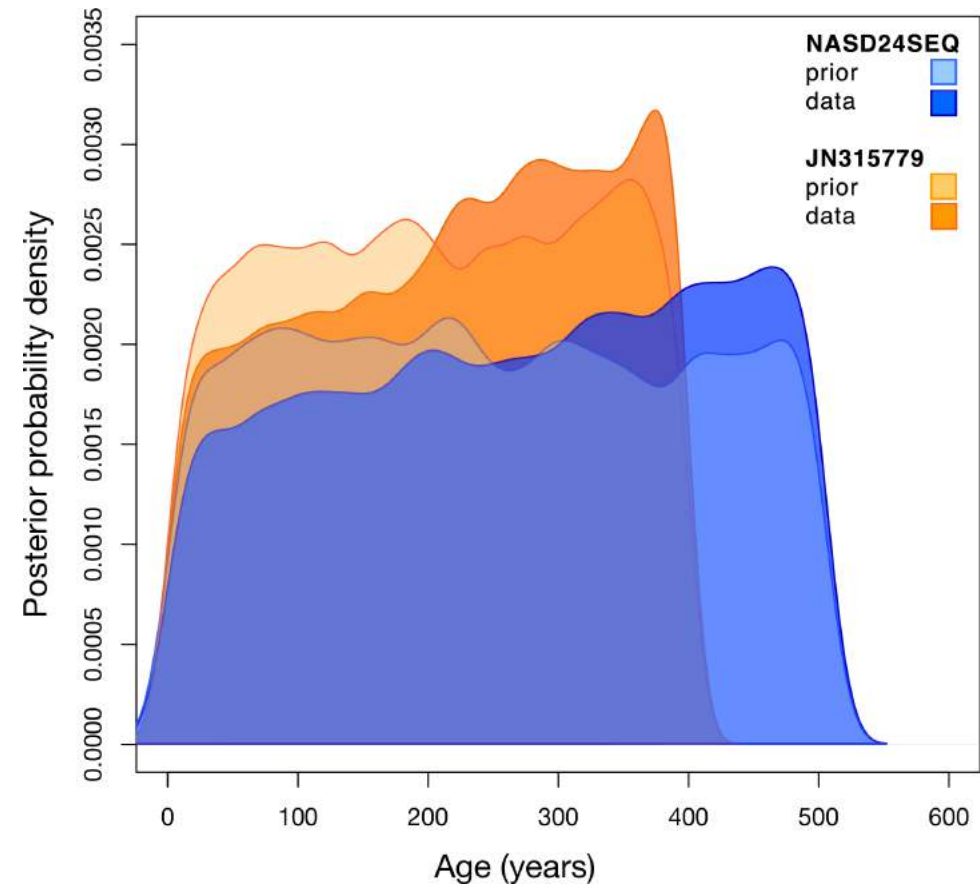
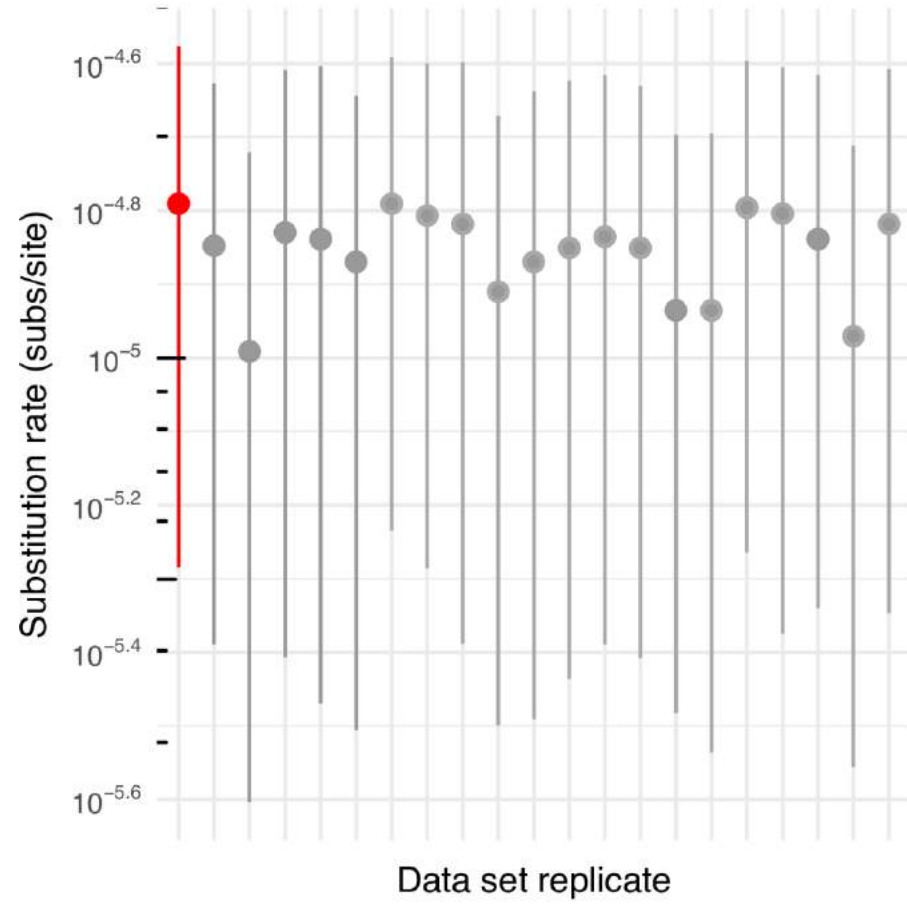


Comparing the prior and posterior -tips



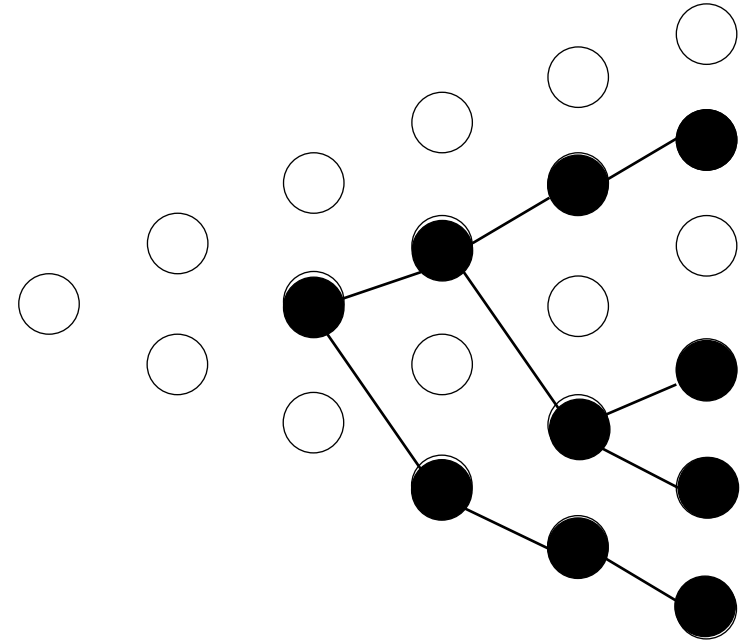
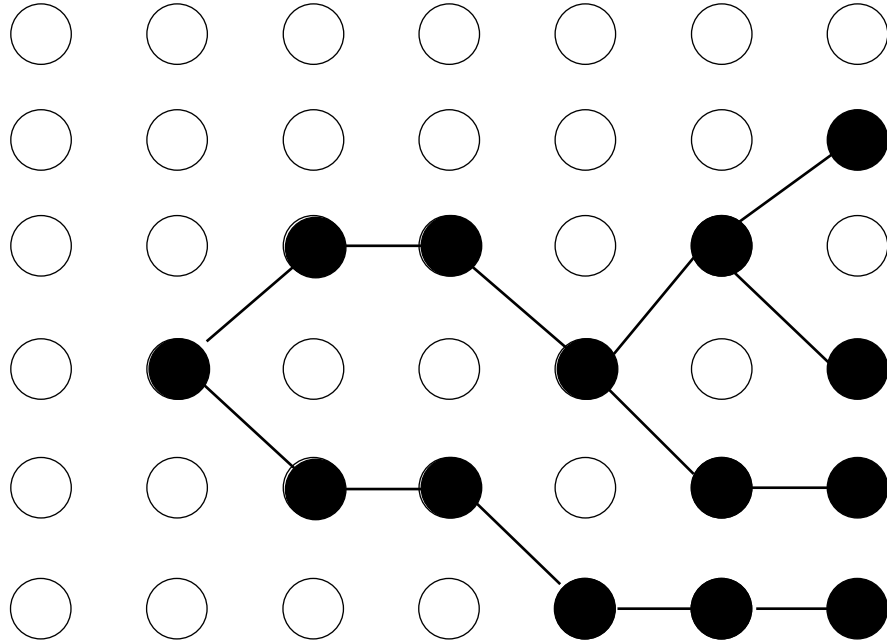
Comparing the prior and posterior -tips





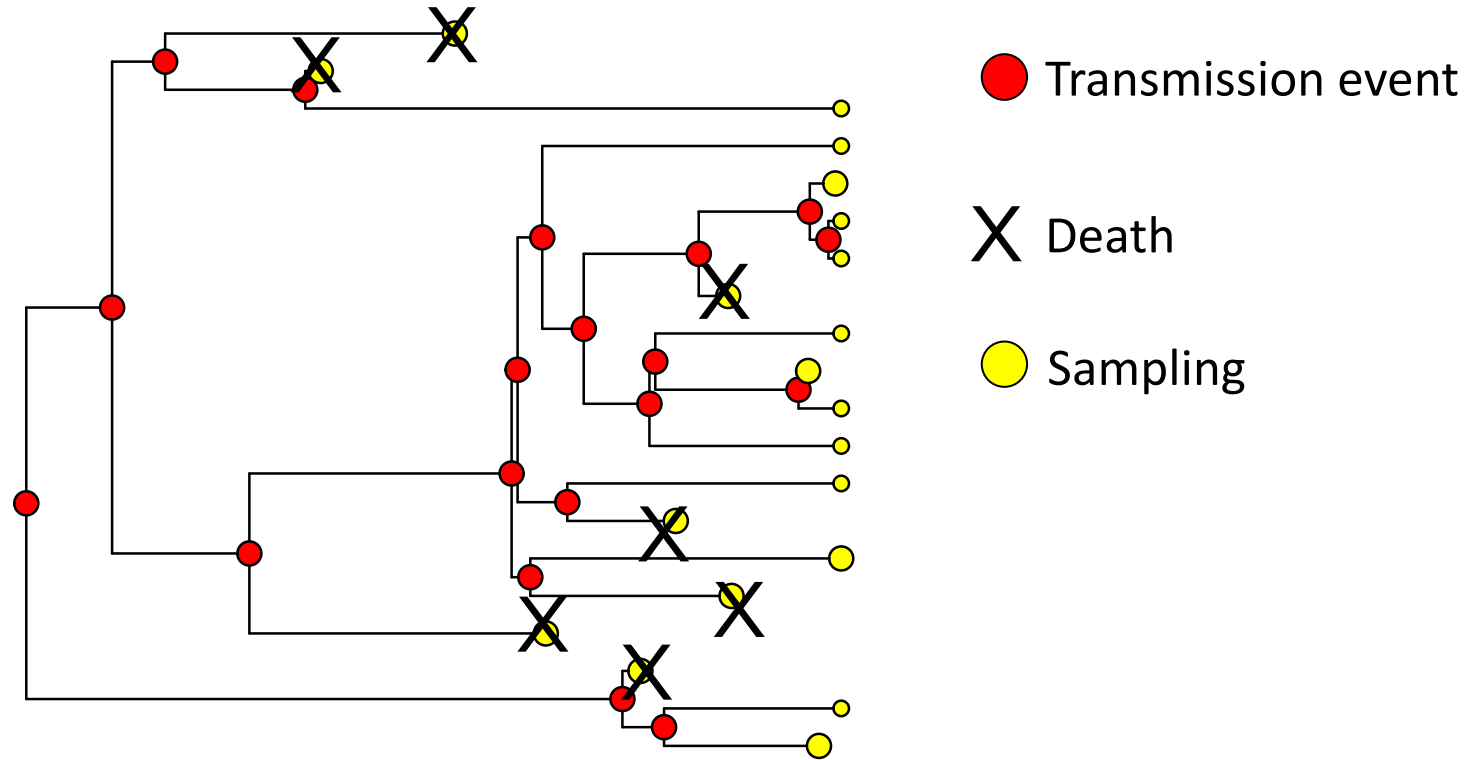
Tree priors and tip-dating

The coalescent



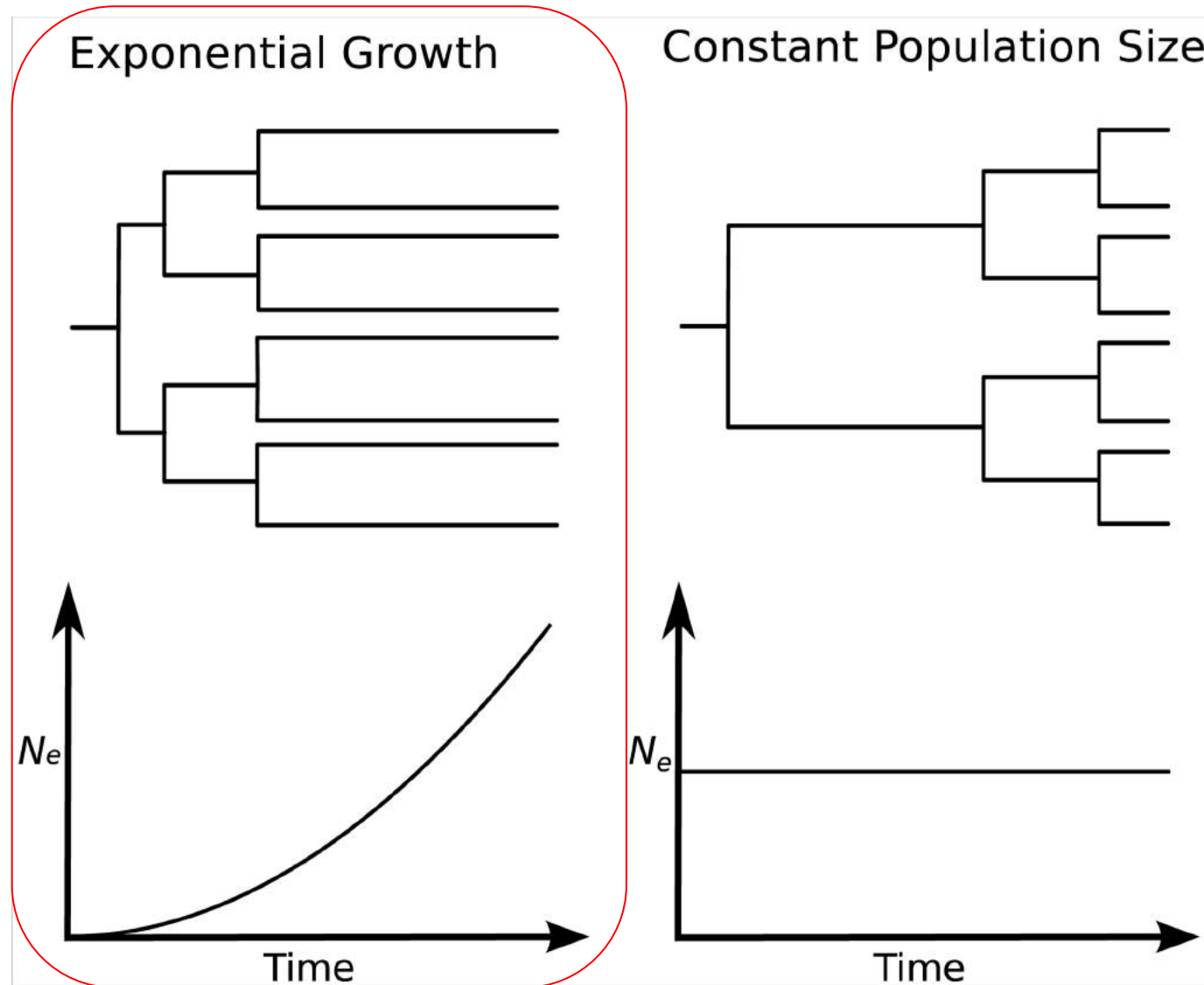
Backwards-in-time process

The birth-death sampling through time



Forward-in-time process

Both, the birth-death with sampling and the exponential coalescent model exponential growth!



Epidemiological parameters

- The coalescent and birth-death have different parameterisations.

Coalescent exponential

r : growth rate

Φ : effective population size

$I(0)$: infected population size
at present

$1 / \delta = \text{duration of infection}$

$$\Phi = I(0) / 2\lambda$$

$$r = \lambda - \delta$$

$$R_0 = (r + \delta) / \delta$$

λ : transmission rate

ψ : sampling rate

μ : death rate

Birth-death

R_0 : **basic reproductive number**

δ : become uninfected rate

ρ : sampling probability

$$R_0 = \lambda / \delta$$

$$\delta = \mu + \psi$$

$$\rho = \psi / (\mu + \psi)$$

Time tree given branching model Prior on branching parameters

$$P(\text{Alignment} \mid \text{Time tree}, \text{Branching model}, \text{Substitution model}, \text{Clock model}) = \frac{P(\text{Alignment} \mid \text{Time tree}) P(\text{Branching model} \mid \text{Substitution model}) P(\text{Substitution model}) P(\text{Clock model})}{P(\text{Alignment})}$$

Alignment Time tree Branching model (can be an epi model) Substitution model Clock model

In the birth-death and the exponential coalescent $P(\text{Time tree} \mid \text{Branching model})$ depends on $\lambda - \delta$ and $\lambda\delta p$ so we usually need a prior on any one parameter to calculate any individual parameter.

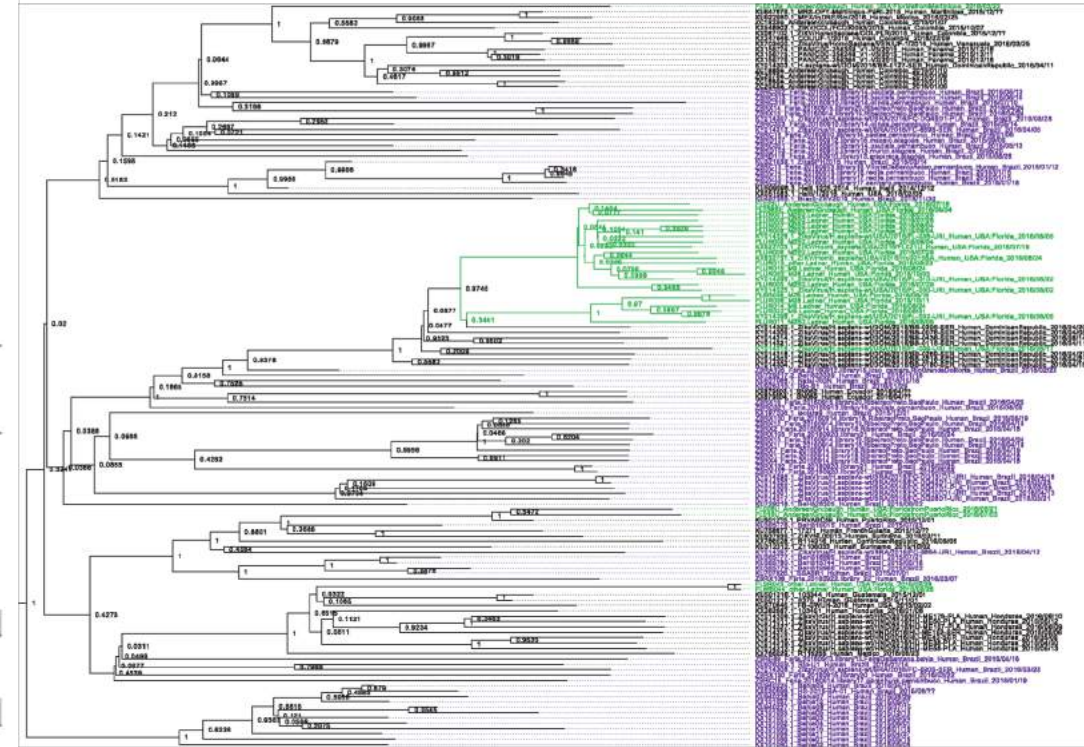
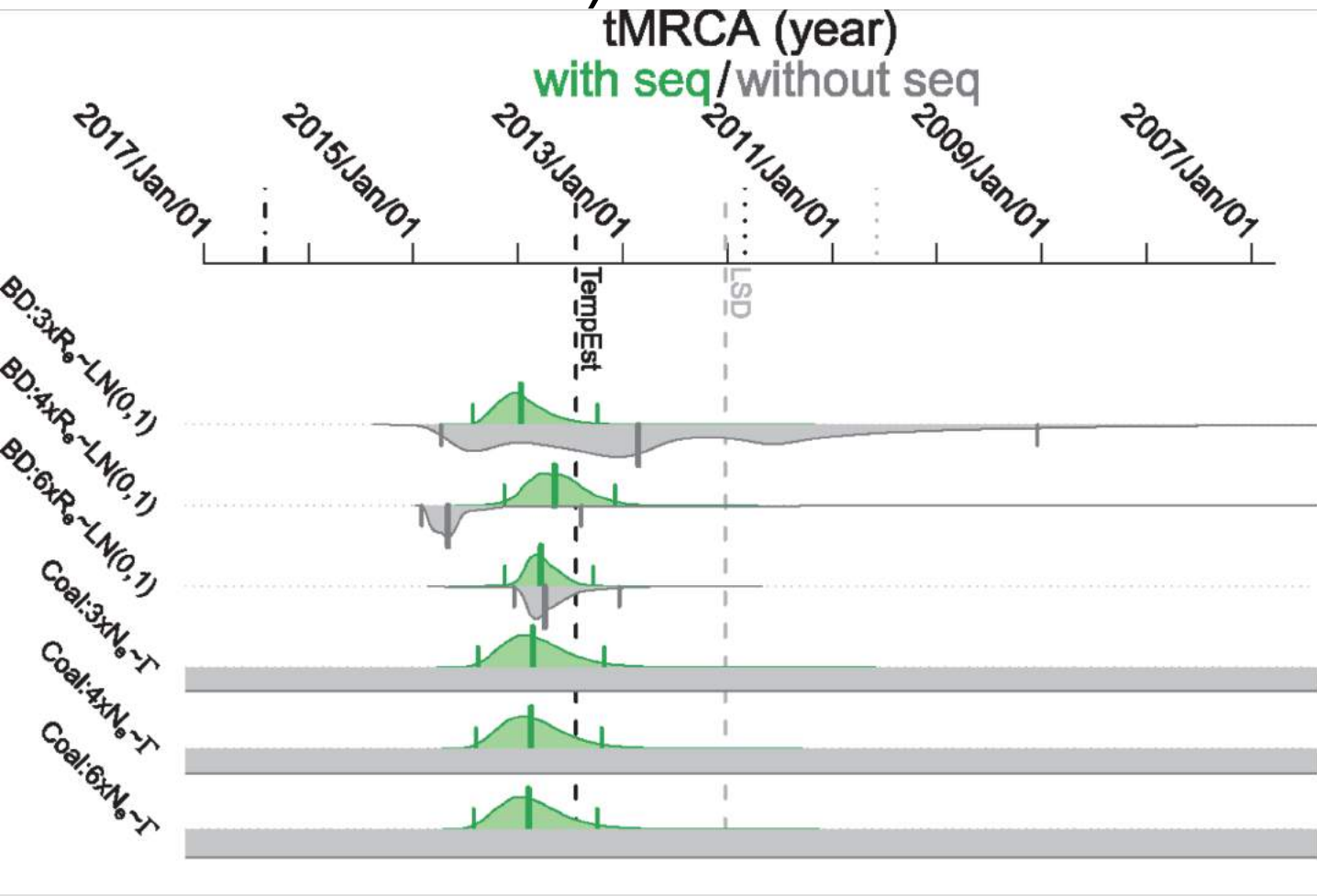
Tree prior $P(\mathcal{E} | \infty)$

$$P_{\text{CoalSamp}}(\mathcal{E} | \text{par} = \Phi, r)$$

$$P_{\text{birth-death}}(\mathcal{E}, \text{sampling times, } N | \text{par} = R_0, \delta, p)$$

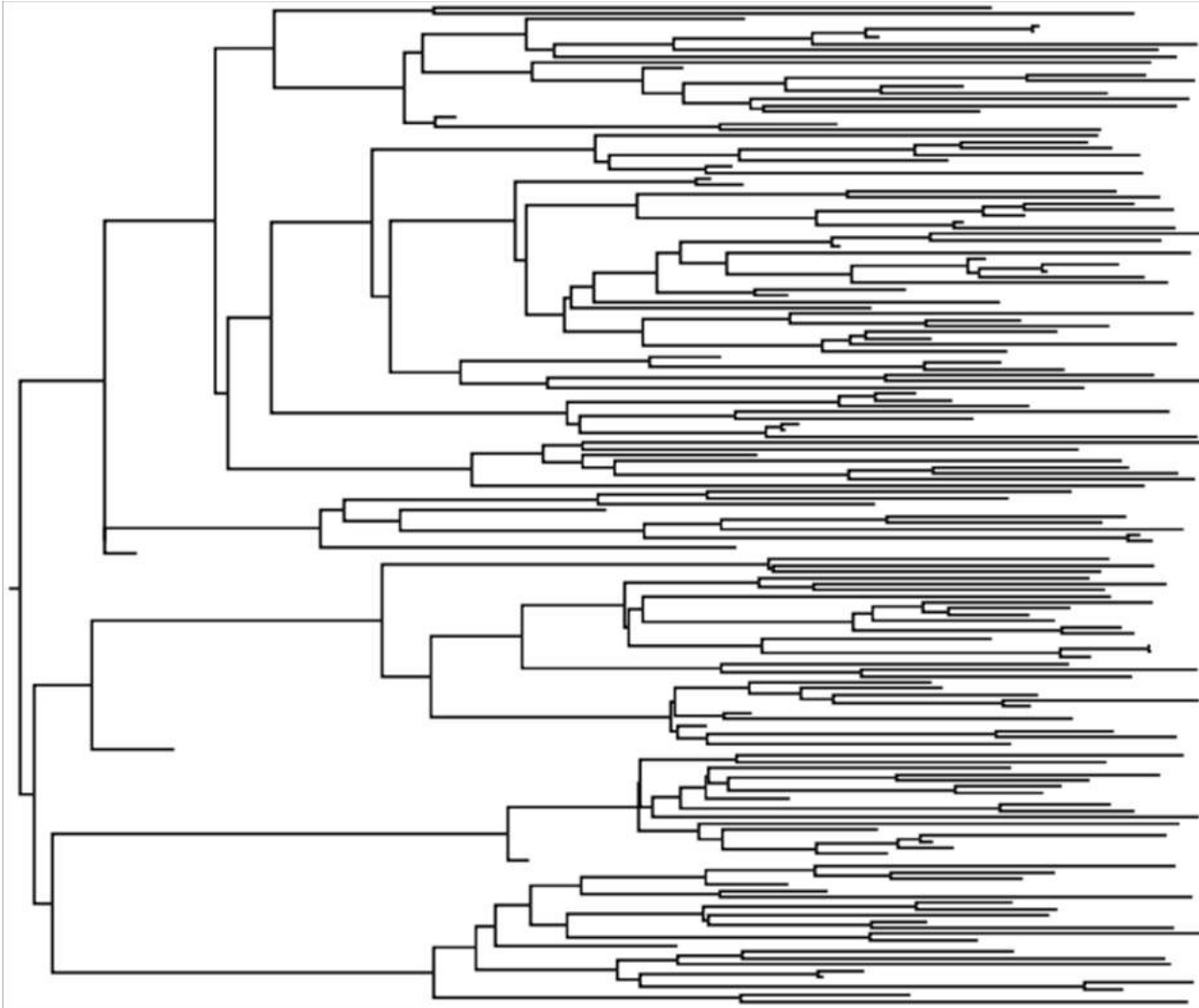
In the birth-death the number of samples and their ages are informative (treated similar to data).

Tree prior and tree height (Zika virus outbreak)



From: Boskova et al. 2018

Tree prior and epi parameters

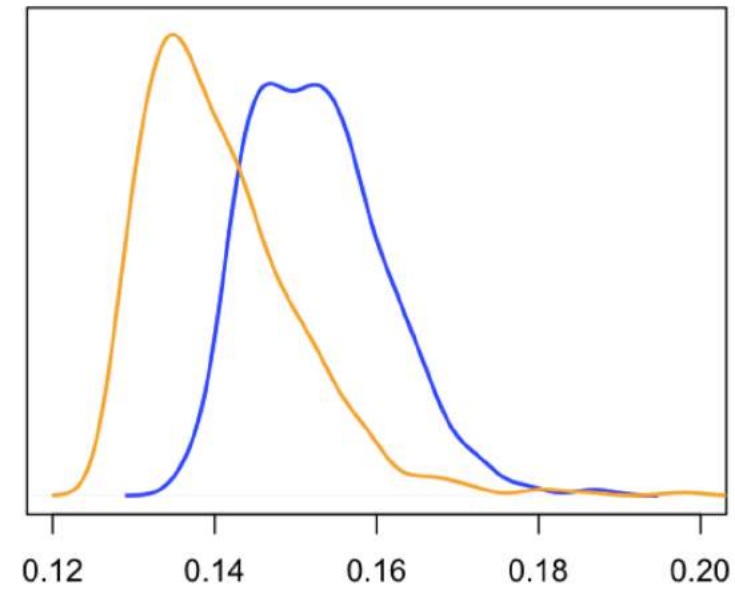
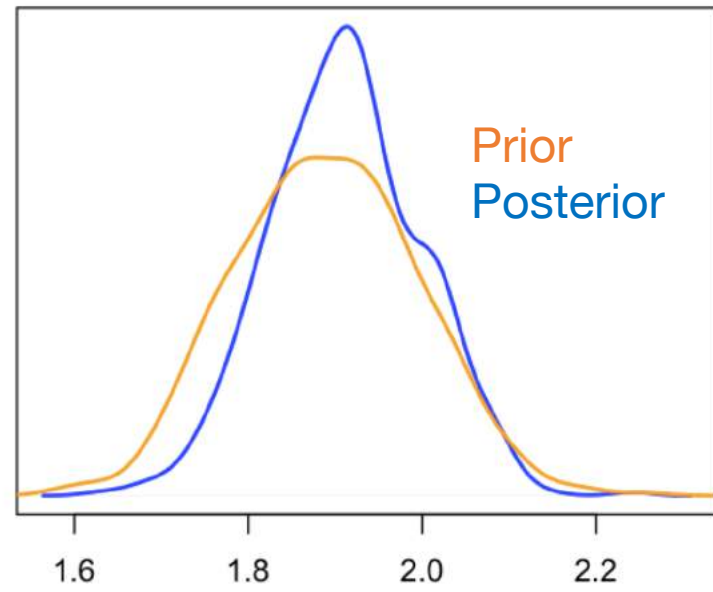


From: di Giallonardo et al. 2018

Data set: Human respiratory syncytial virus (RSVA)

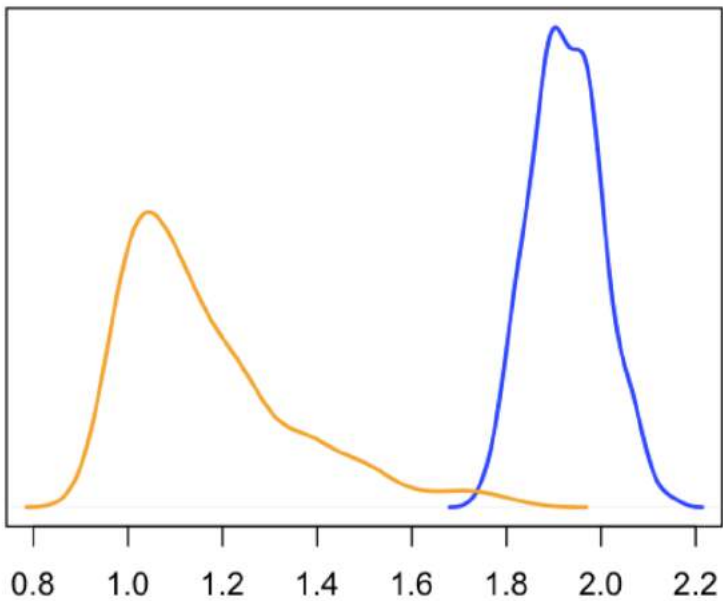
- 163 samples collected over one month.
- Sampled in August (peak number of cases).
- Analyse using birth-death and coalescent exponential.

Birth-death

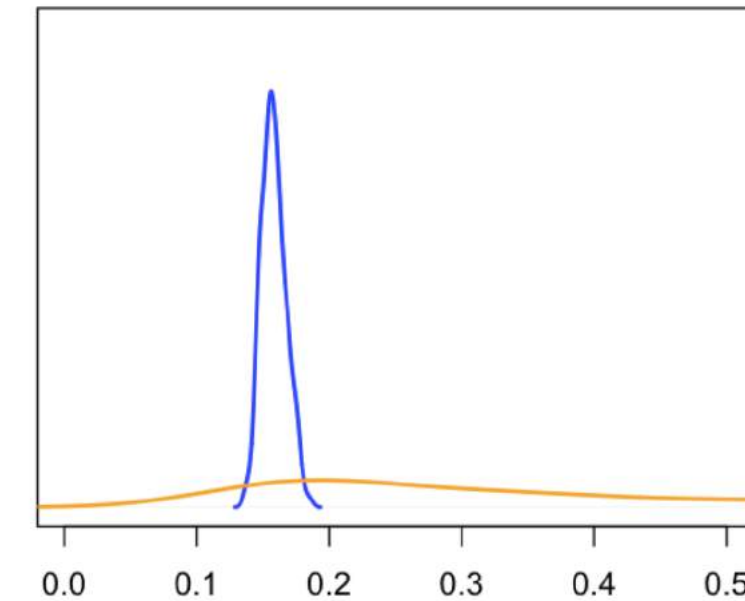


In the birth-death, sampling times alone are very informative about epi parameters.

Coalescent exponential

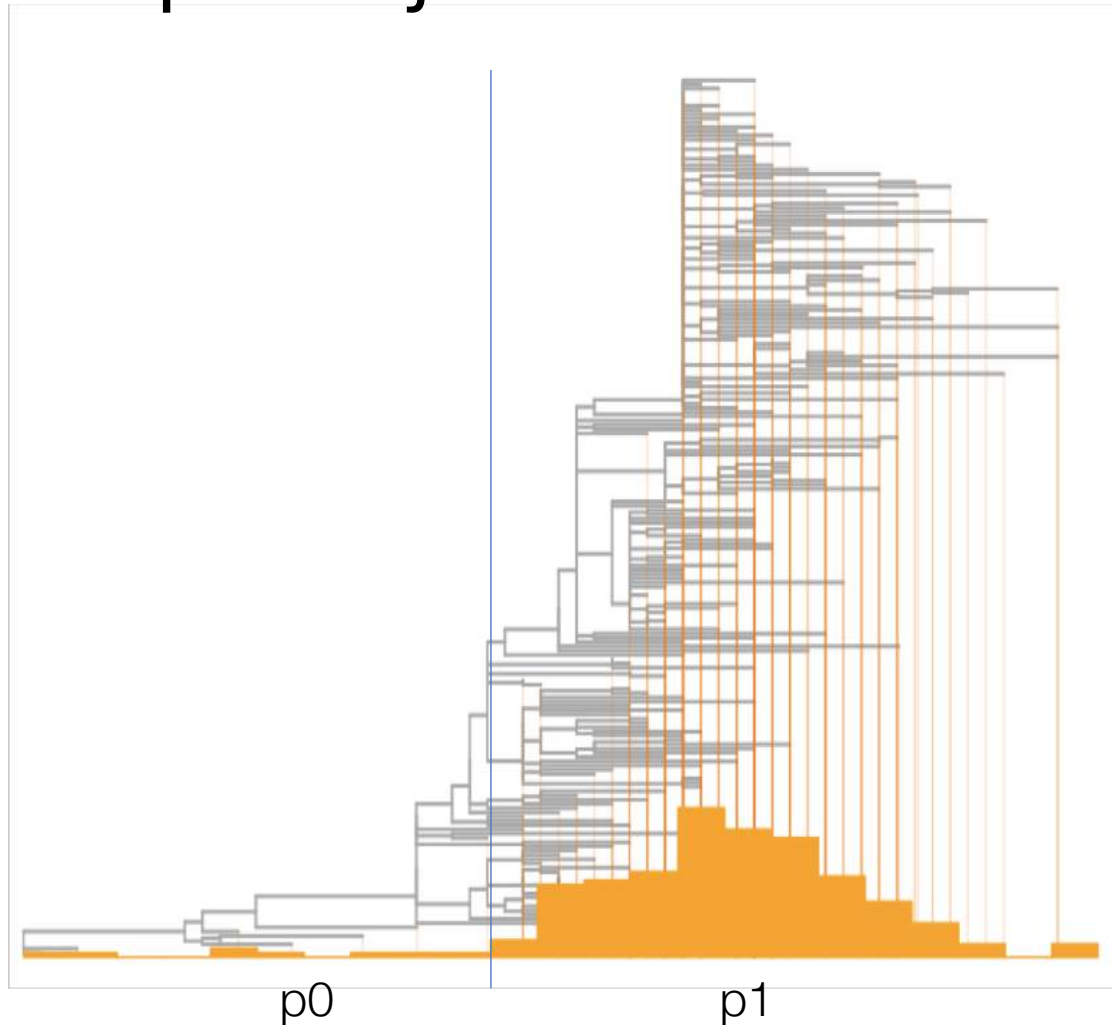


R_0



Root height (years)

Epi trajectories



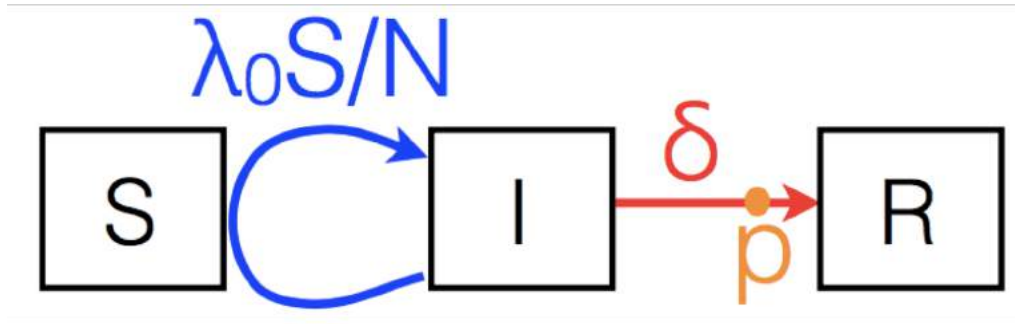
Account for sampling heterogeneity using BDSIR (or BD skyline).

Data set: AH1N1 Influenza

- Samples collected in 2009 from February to August
- Estimate R_0 using three model:
 - Birth-death
 - Coalescent exp.
 - Birth-death SIR

Susceptible-infected-recovered

$$N=S+I+R$$

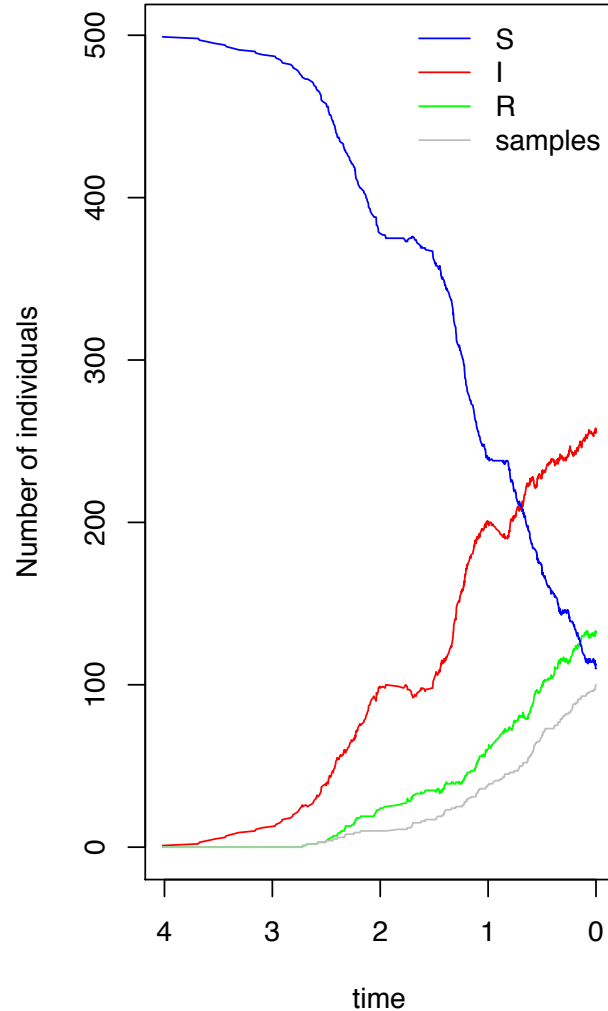


λ = transmission rate

δ = becoming non-infectious rate

ρ = sampling probability

Susceptible-infected-recovered



Transmission rate (birth): β

Become uninfected rate: γ

Susceptible pop. size: n_s

$$R_e = \beta * n_s / \gamma$$

$$R_0 = \beta * n_s(0) / \gamma$$

SIR and BD are similar when
 $S \rightarrow N$

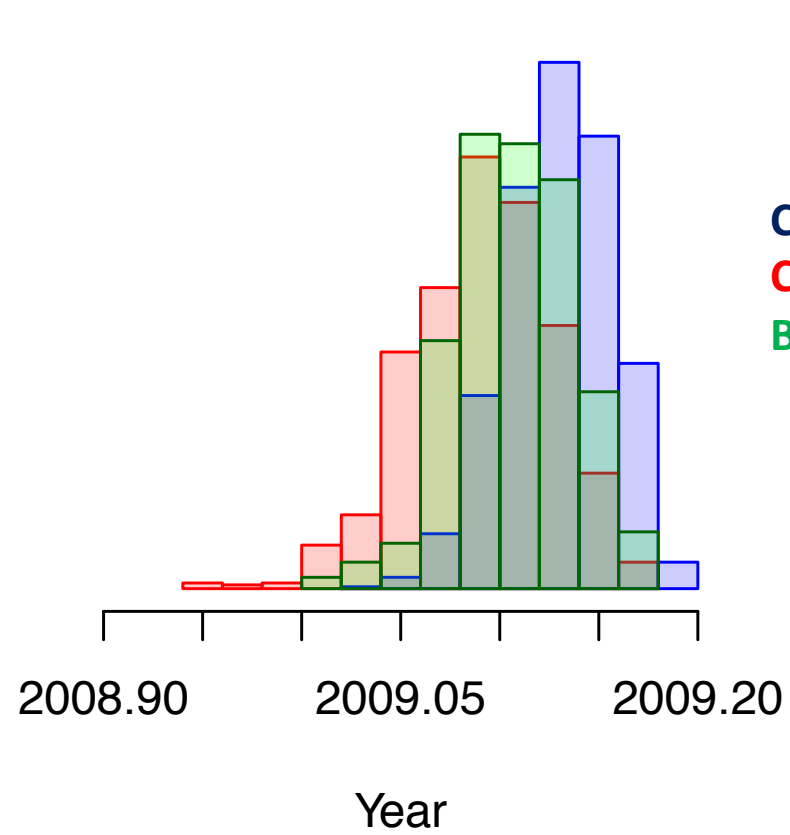
Comparing to birth-death notation:

γ is equivalent to δ

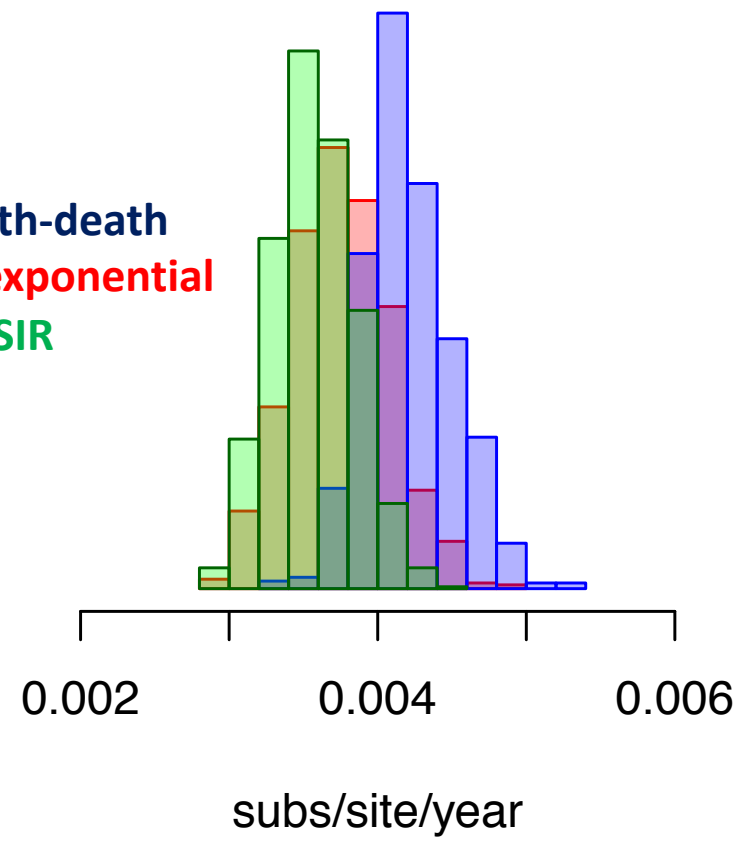
$$\beta * n_s(0) = \lambda_0$$

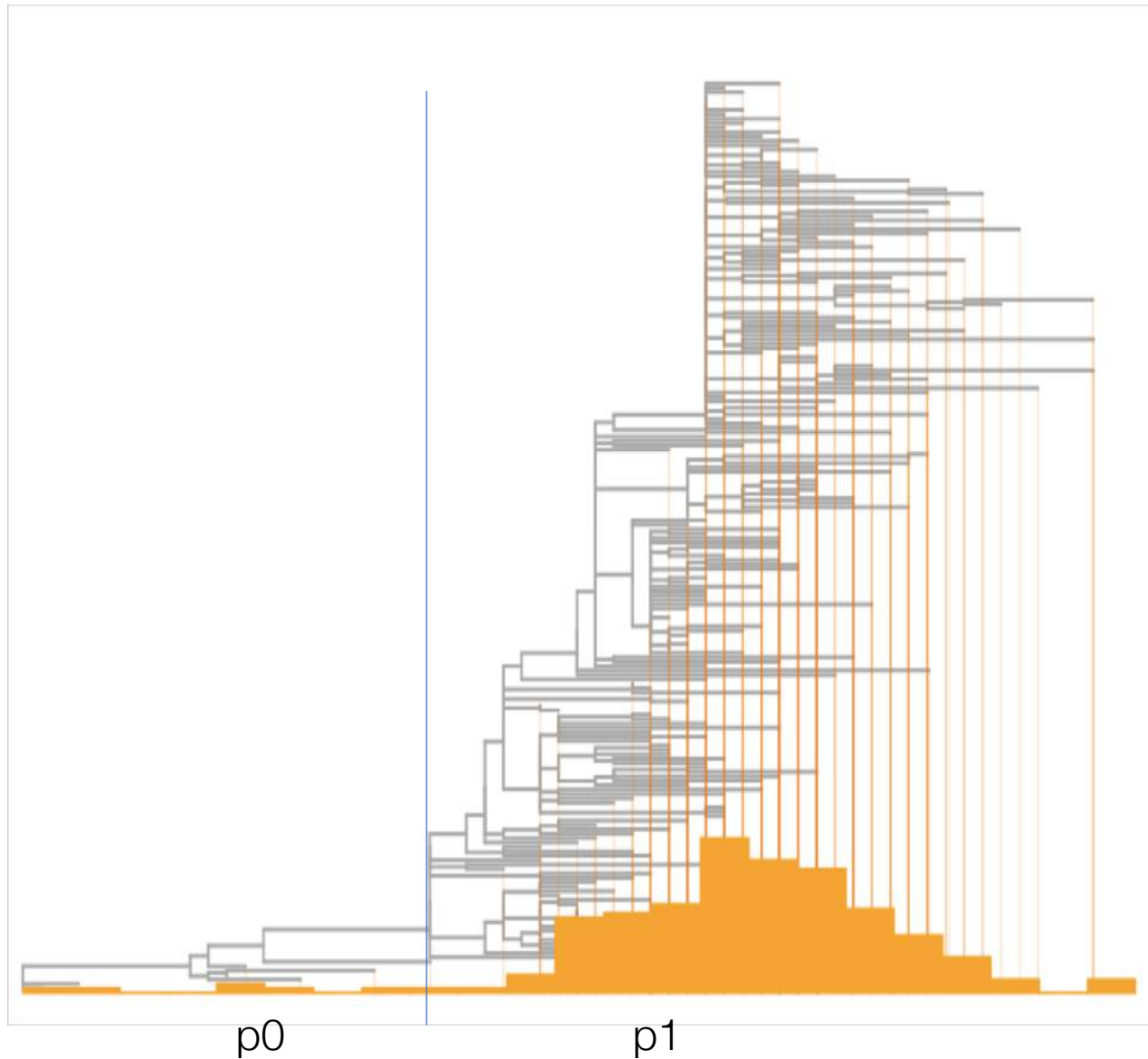
The birth-death SIR is an approximation of the SIR based on the birth-death skyline (also available; coalSIR, SIR).

Age of samples

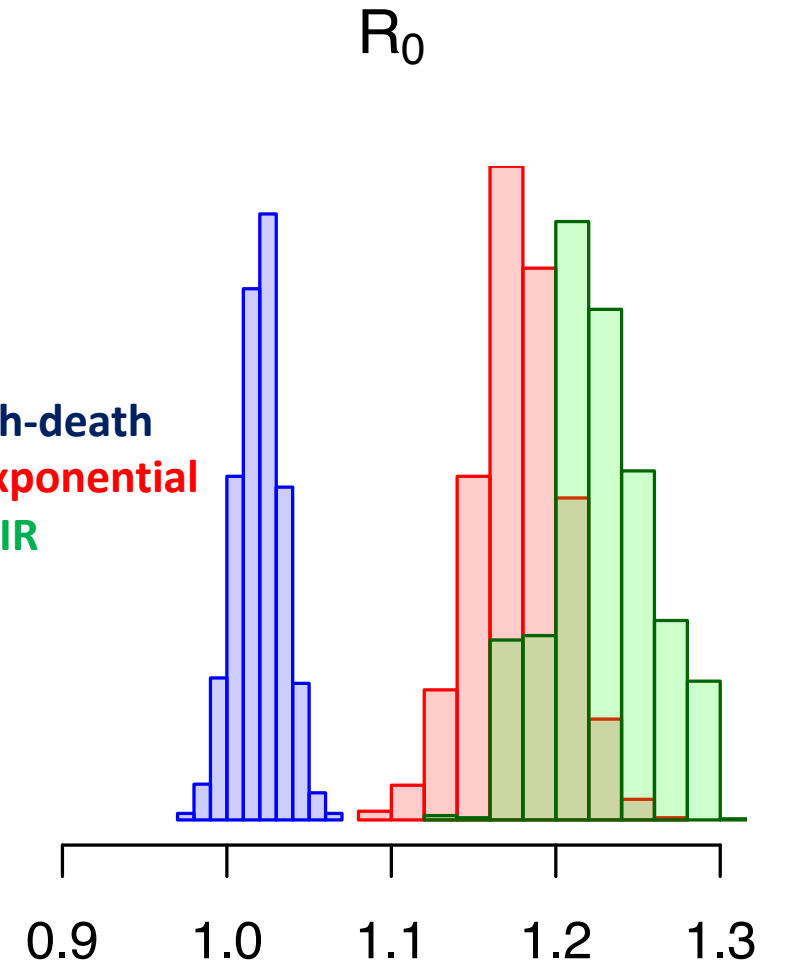


Evol. rate



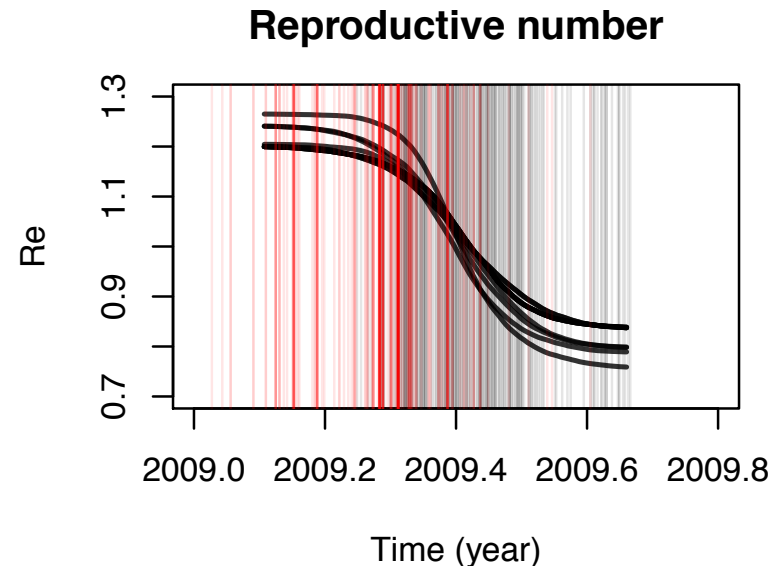
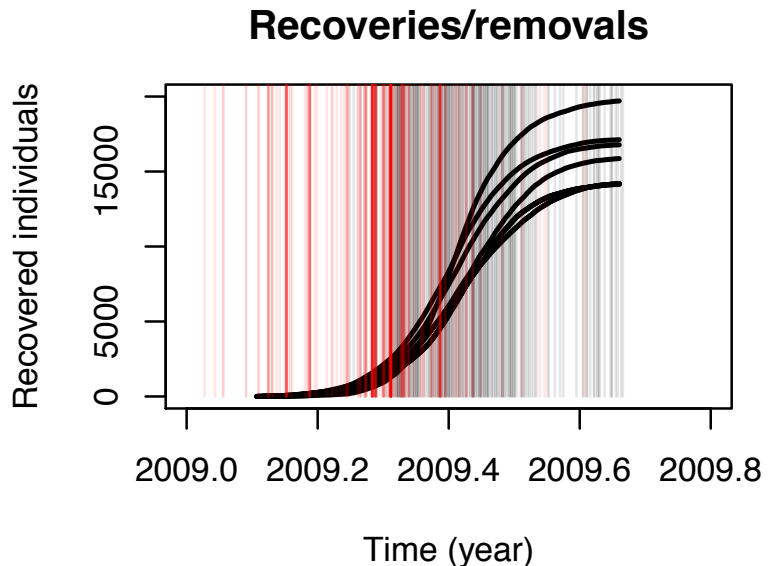
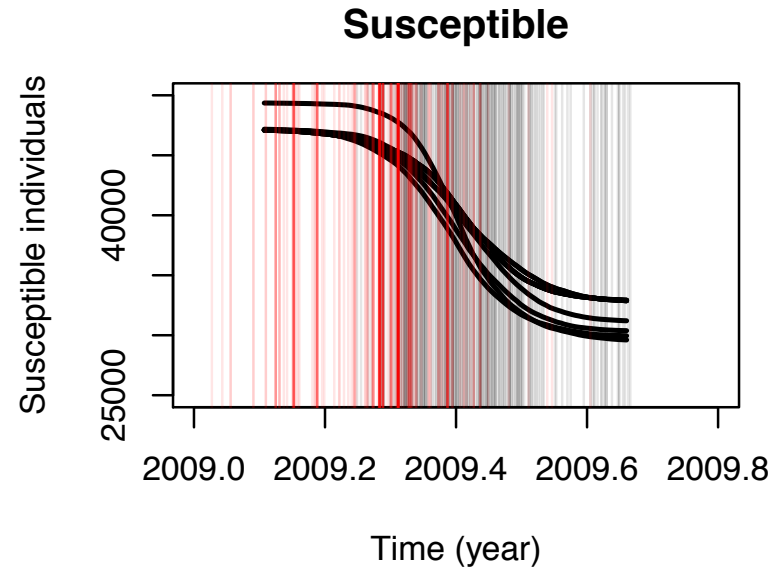
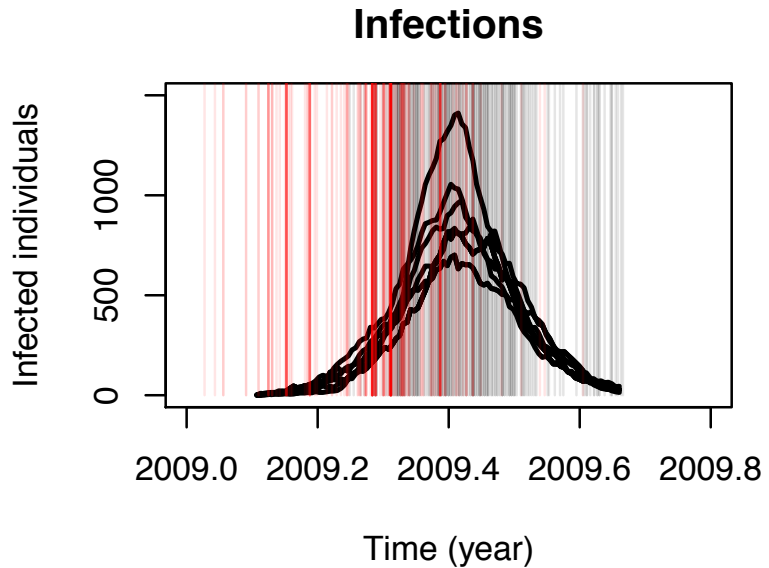


Constant Birth-death
 Coalescent exponential
 Birth-death SIR



The **constant birth-death** assumes constant sampling over time!
 We can relax this assumption using the birth-death Skyline or the
 birth-death SIR.

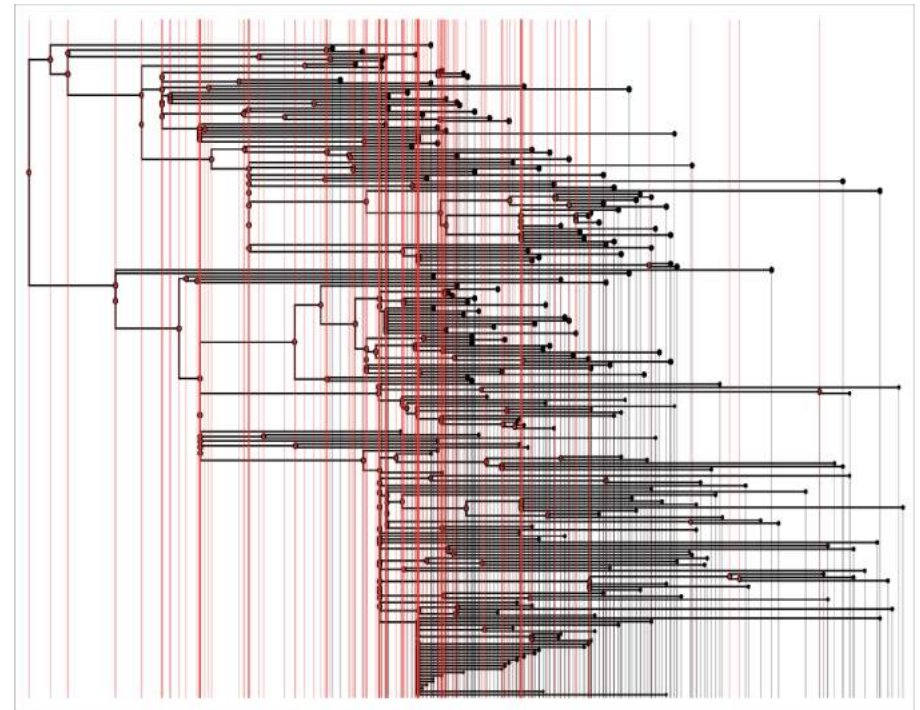
SIR trajectories (BDSIR, CoalSIR, or SIR)



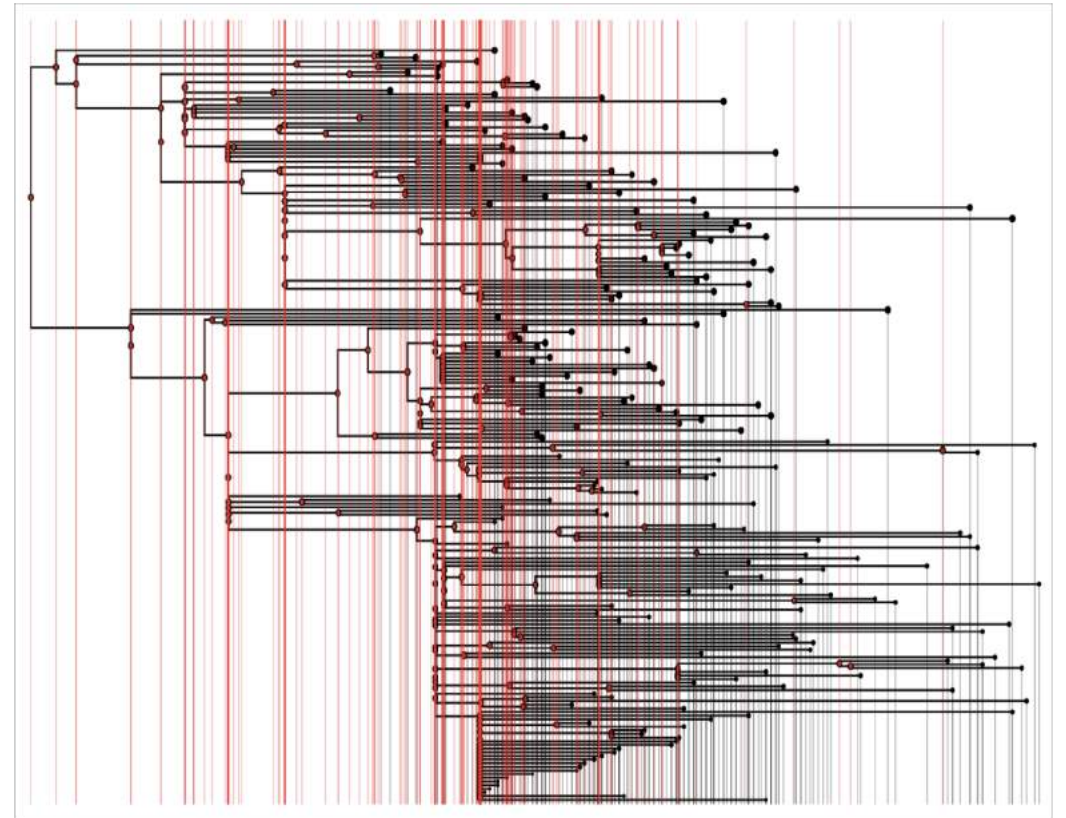
- Collected genome samples

- Internal nodes
(inferred transmission events)

- Tip-dates are a rich source of information
 - Calibration.
 - Information for some tree priors (e.g. birth-death).
- Verify temporal structure
 - Comparing prior and posteriors.
 - Tip-dates can be specified using probability distributions (e.g. ancient DNA).



- Epidemiological inferences
 - Internal node information.
 - Branching as an epidemiological process
 - Check phylodynamic model assumptions
 - E.g. via model adequacy; phyloseminar.org on model adequacy of infectious disease phylodynamics and Duchene et al. 2018.
- In the future...
 - Hospital notification data.
 - Drug resistance profiles.
 - Virulence.



Some useful references

- Boskova, V., Stadler, T., & Magnus, C. (2018). The influence of phylodynamic model specifications on parameter estimates of the zika virus epidemic. *Virus evolution*, 4(1), vex044.
- Du plessis, I., & Stadler, T. (2015). Getting to the root of epidemic spread with phylodynamic analysis of genomic data. *Trends in microbiology*, 23(7), 383-386.
- Ho, s. Y., & Duchêne, S. (2014). Molecular-clock methods for estimating evolutionary rates and timescales. *Molecular ecology*, 23(24), 5947-5965.
- Duchêne, s., Duchêne, D., Holmes, E. C., & Ho, S. Y. (2015). The performance of the date-randomization test in phylogenetic analyses of time-structured virus data. *Molecular biology and evolution*, 32(7), 1895-1906.

Model adequacy and phylogenetics

Model selection

M1, M2, M3

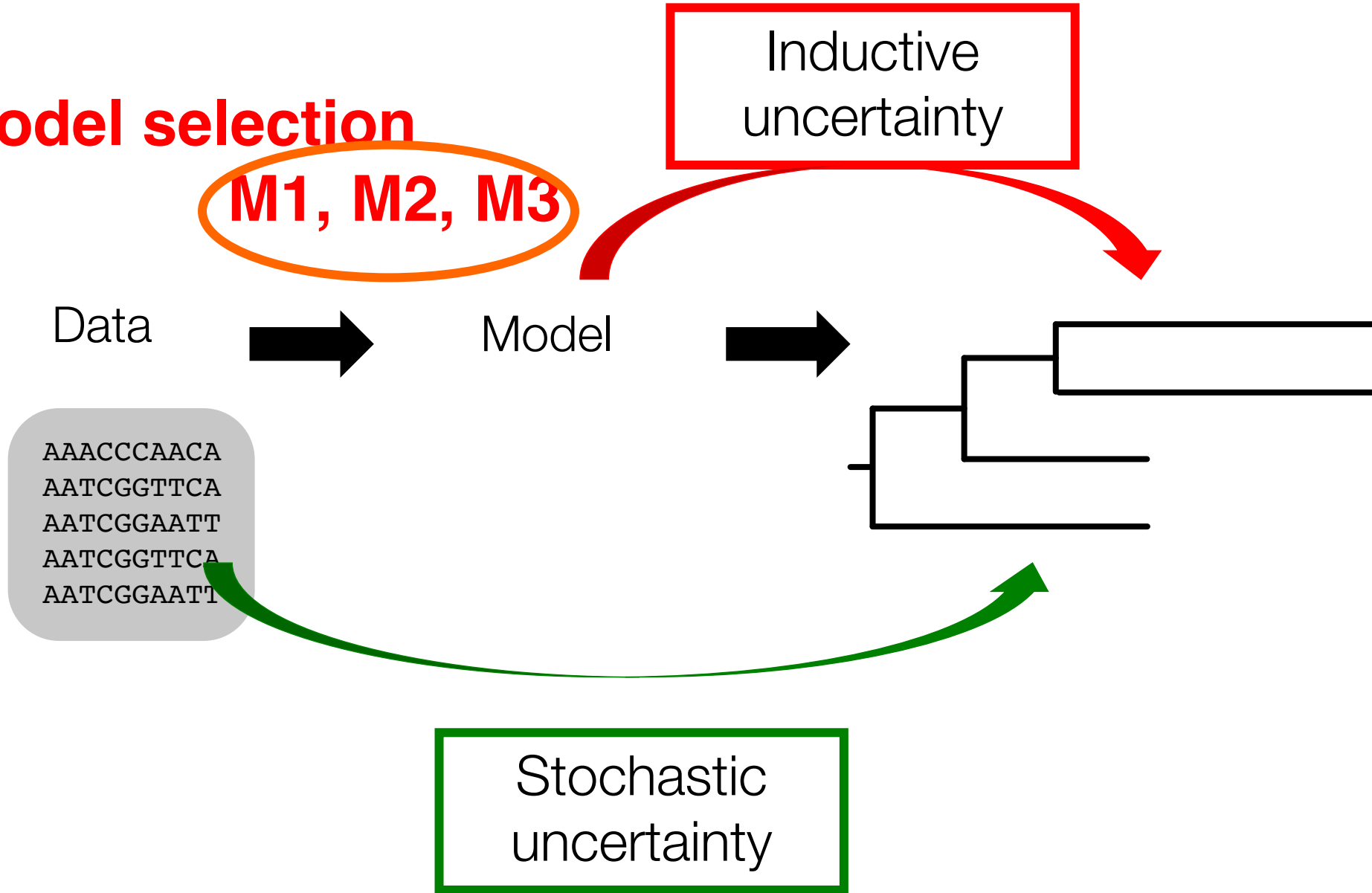
Inductive
uncertainty

Data

Model

AAACCCAACA
AATCGGTTCA
AATCGGAATT
AATCGGTTCA
AATCGGAATT

Stochastic
uncertainty



Model adequacy

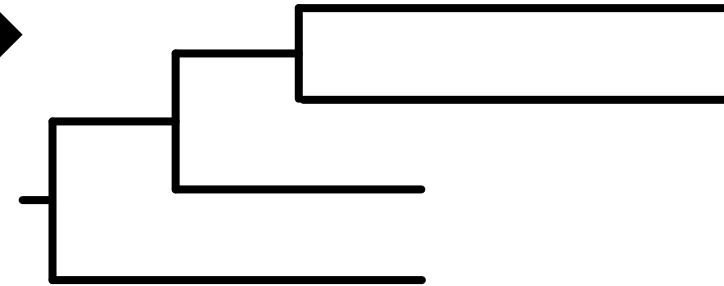
~~M1~~, ~~M2~~, ~~M3~~

Inductive
uncertainty

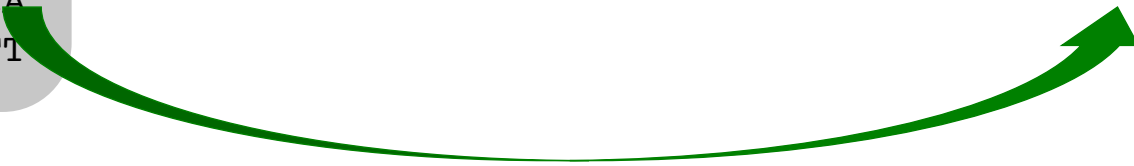
Data



Model



AAACCCAACA
AATCGGTTCA
AATCGGAATT
AATCGGTTCA
AATCGGAATT



Stochastic
uncertainty

Model selection

Select a pool of models

Rank models according to statistical fit

Select that with the highest fit

Model adequacy

Consider a model(s)

Treat models as hypotheses

Conduct hypothesis testing to assess the limits of the model's applicability

The goal of model adequacy

- *The relevant goal is not to answer the question, ‘Do the data come from the assumed model?’ (to which the answer is almost always no), but to **quantify the discrepancies between data and model**, and assess whether they could have arisen by chance, under the model’s own assumptions.*

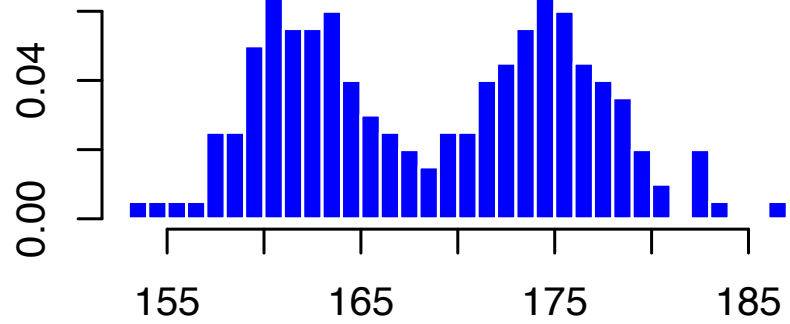
Gelman et al. 2013

Chapter 6 of Bayesian Data Analysis, 3rd Edition



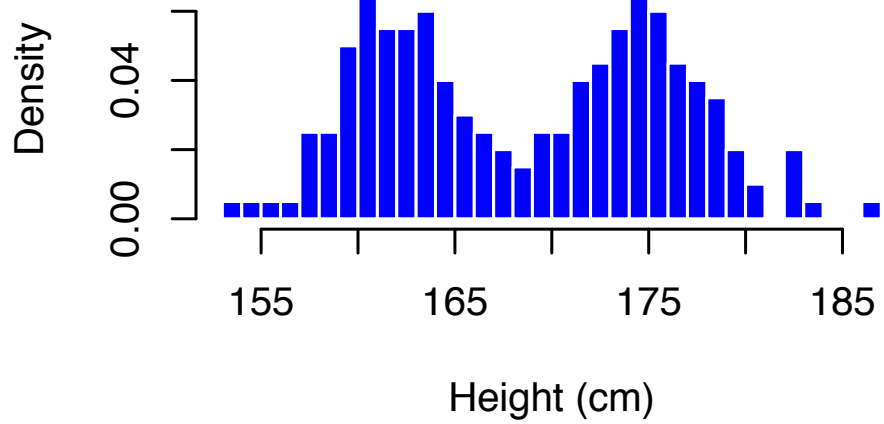
Height

Density

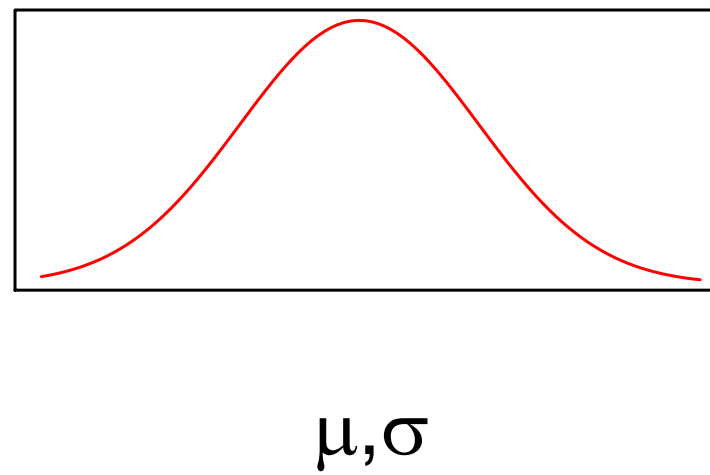


Height (cm)

Height

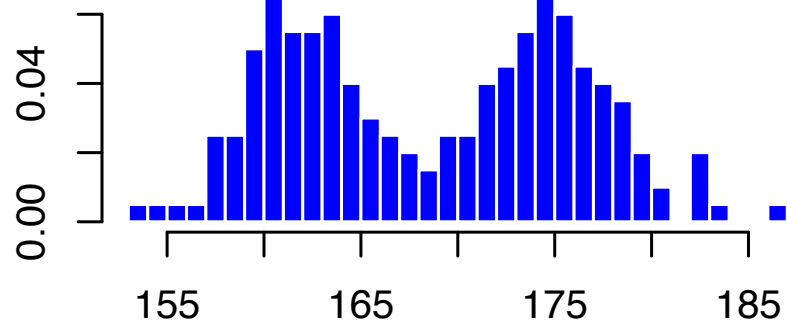


Normal model



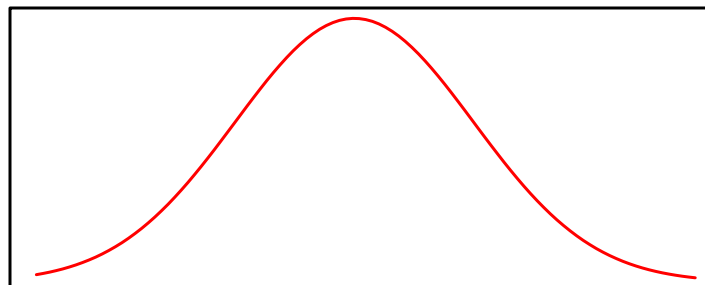
Height

Density

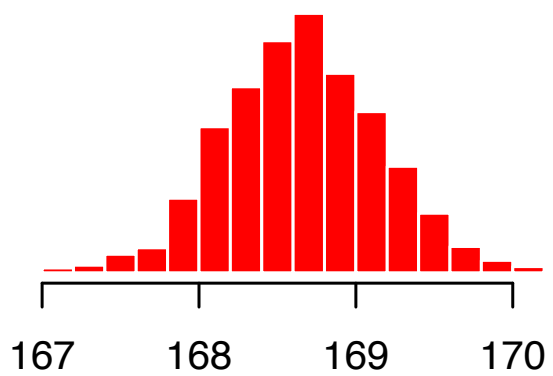


Height (cm)

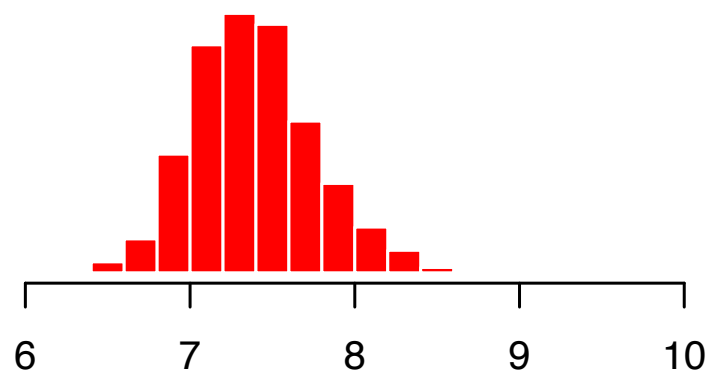
Normal model



μ, σ



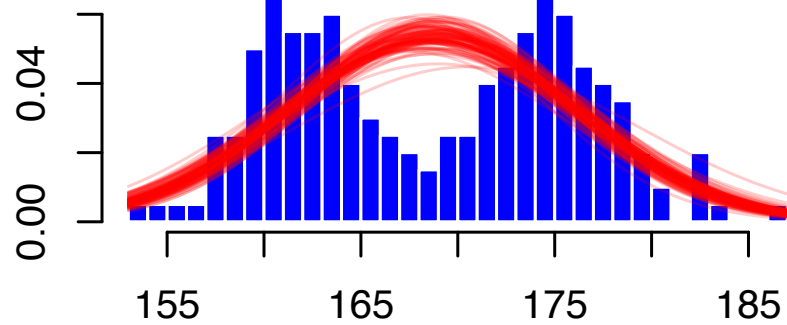
μ



σ

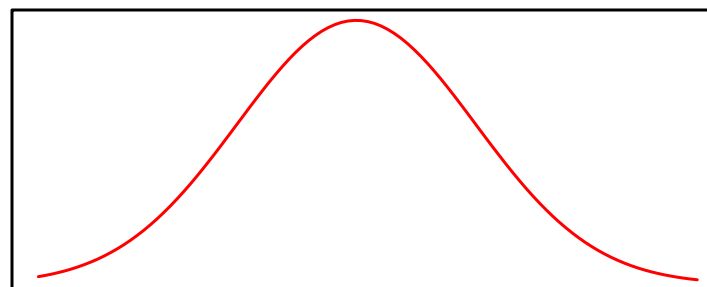
Height

Density

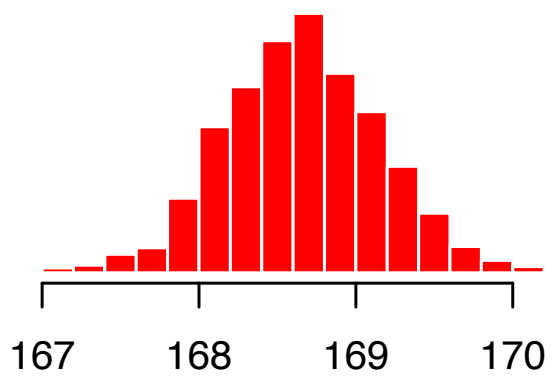


Height (cm)

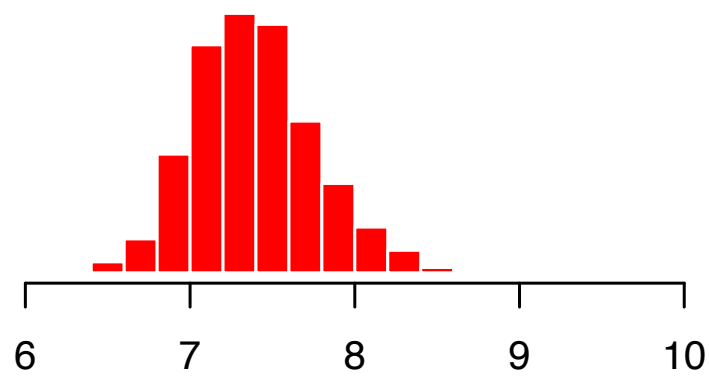
Normal model



μ, σ

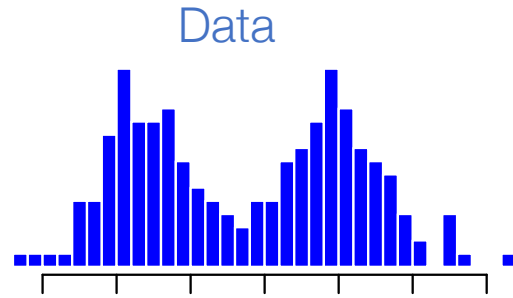


μ

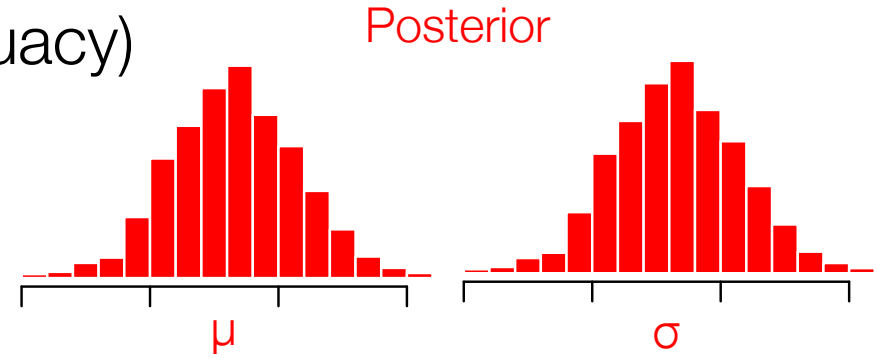


σ

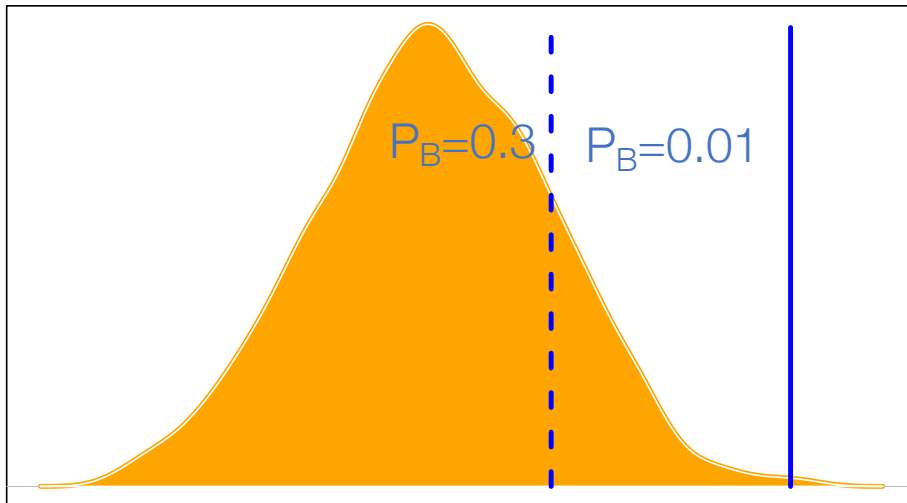
Posterior predictive model checking (Bayesian model adequacy)



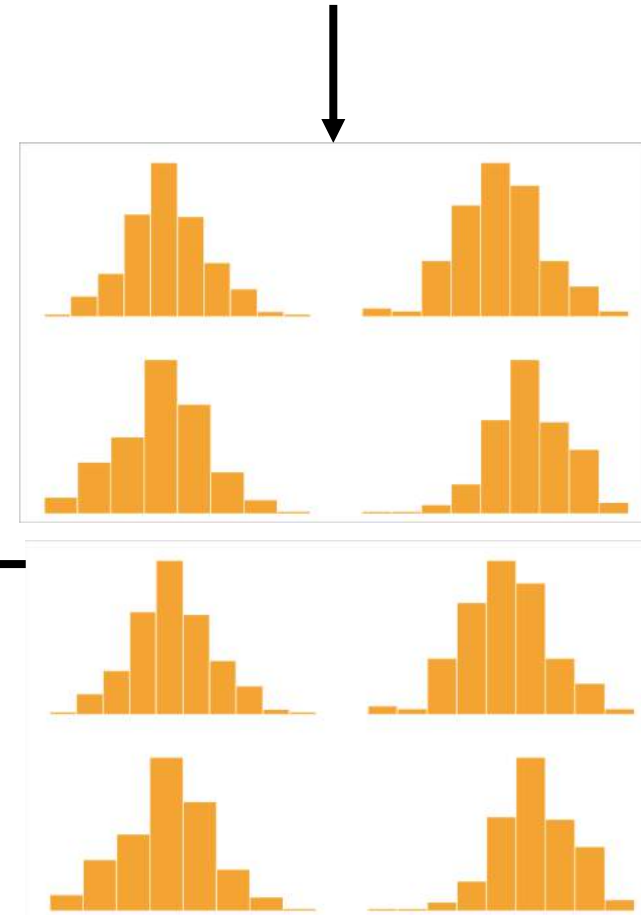
Model M1
 (μ, σ)



Test
statistic



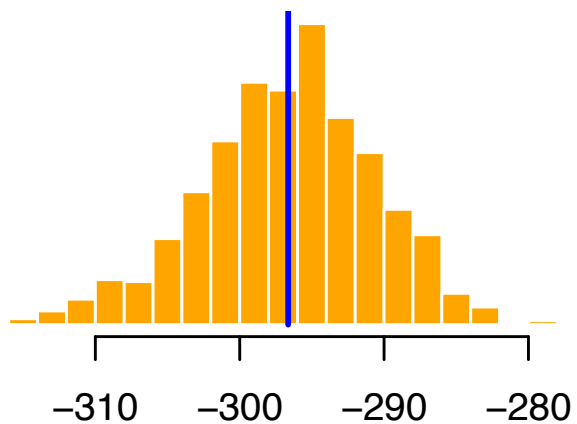
Posterior predictive test



Posterior
predictive
simulations

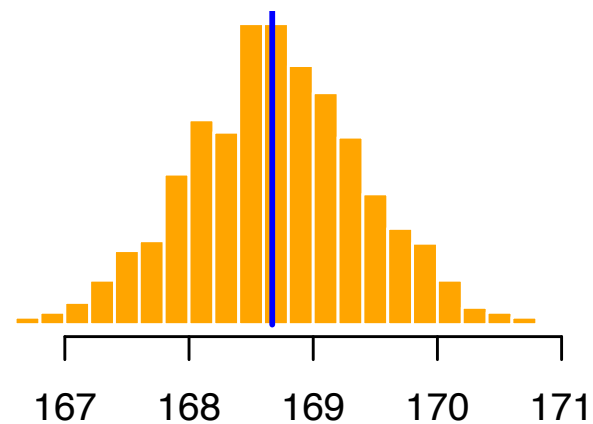
Test
statistic

Mean likelihood



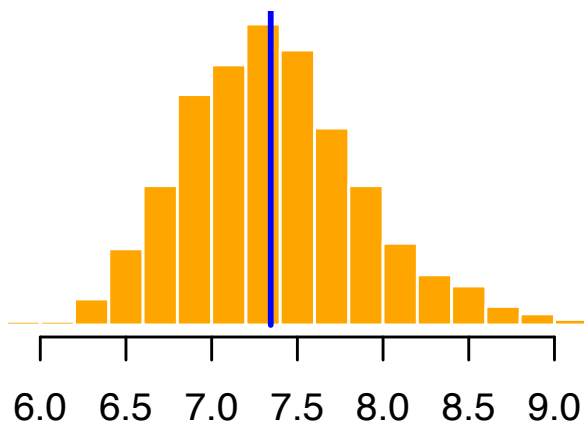
$P_B = 0.481$

Mean



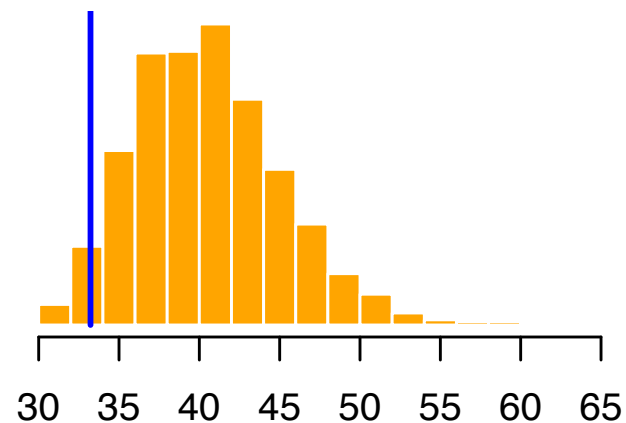
$P_B = 0.494$

Sd



$P_B = 0.522$

Range



$P_B = 0.019$

“Sufficient statistics are bad test statistics”

On multiple comparisons...

- We can account for multiple comparisons (e.g. Bonferroni correction or multidimensional p-value).
- However, we do not make this adjustments because the goal of model adequacy is to assess how the model predicts particular aspects of the data.

Tree prior/ branching model /phylodynamic model

Duchene S et al. 2018, Revell et al. 2005,

Drummond and Suchard 2008

Höhna 2015, Pennell et al. 2015

E.g.

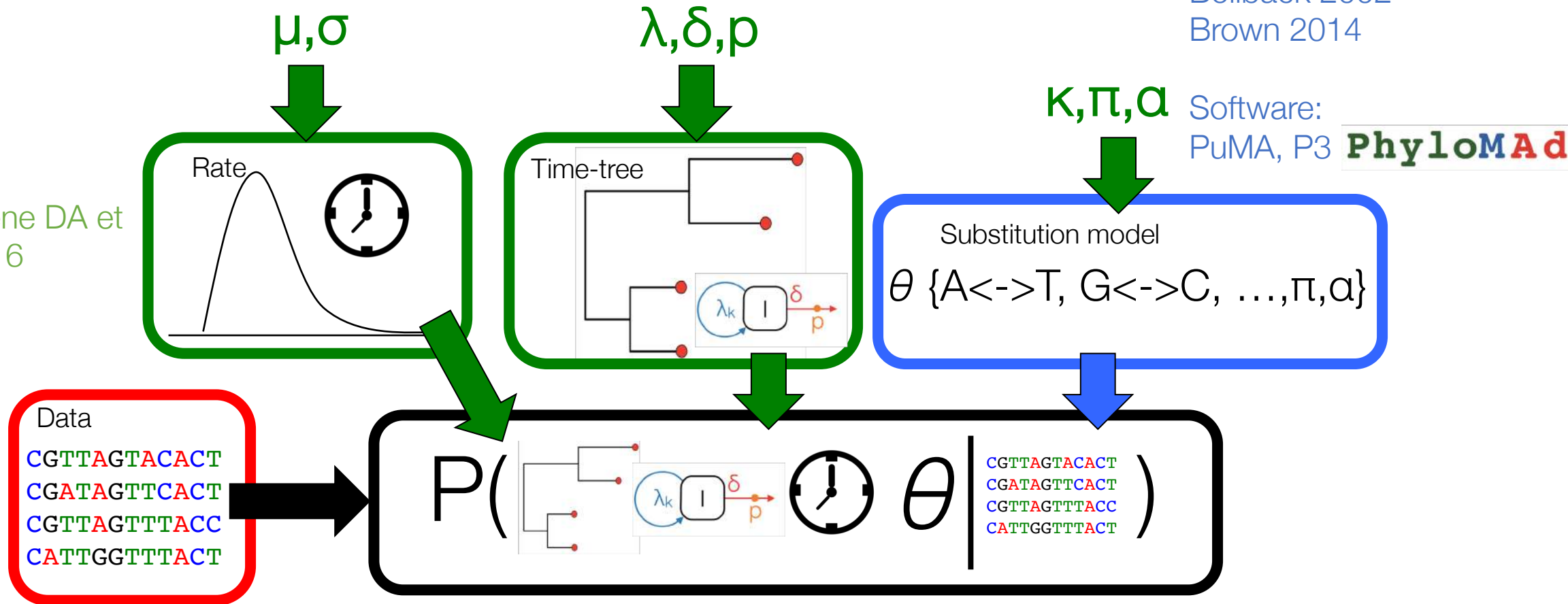
Reeves 1992

Goldman 1993

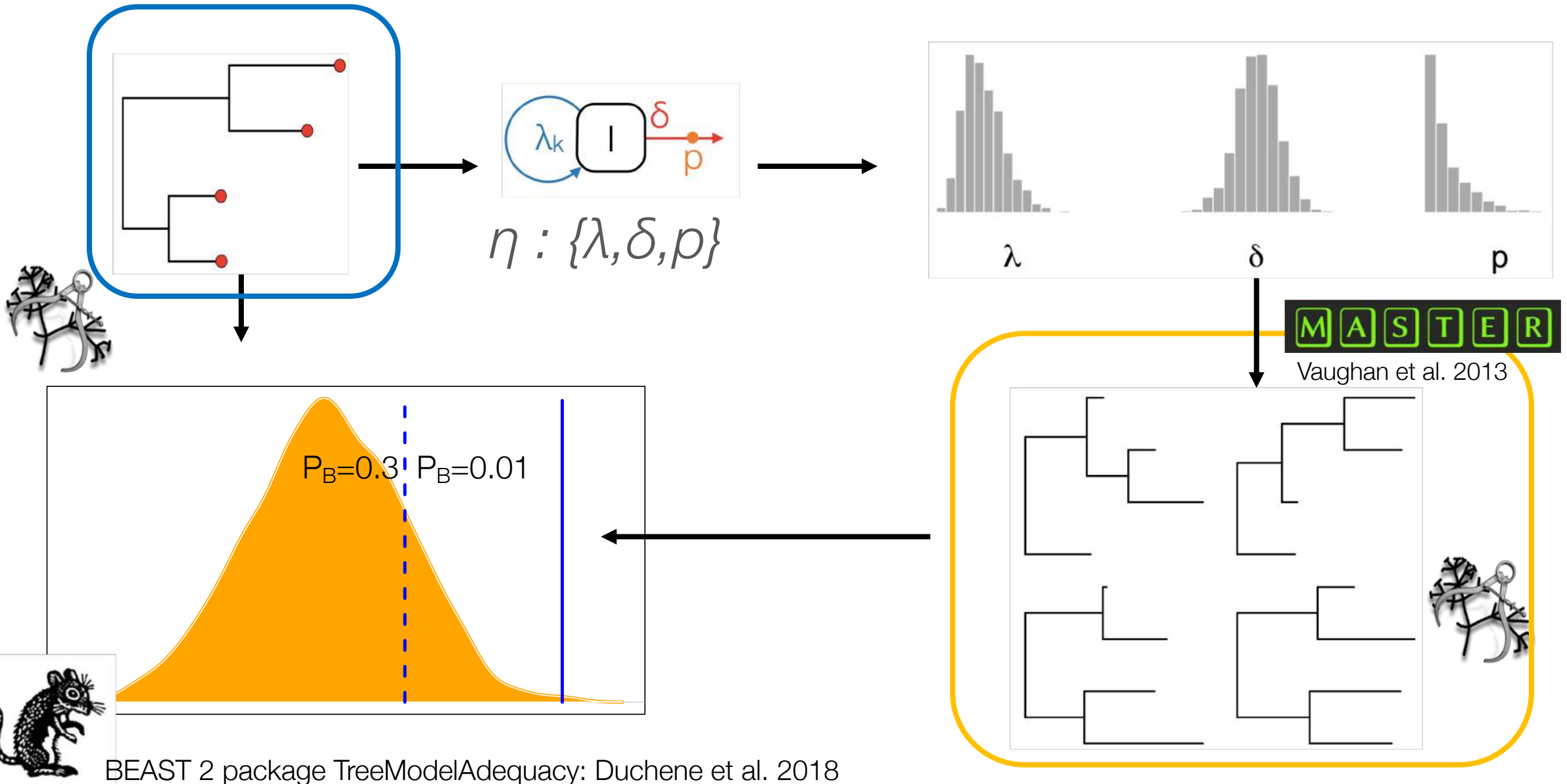
Bollback 2002

Brown 2014

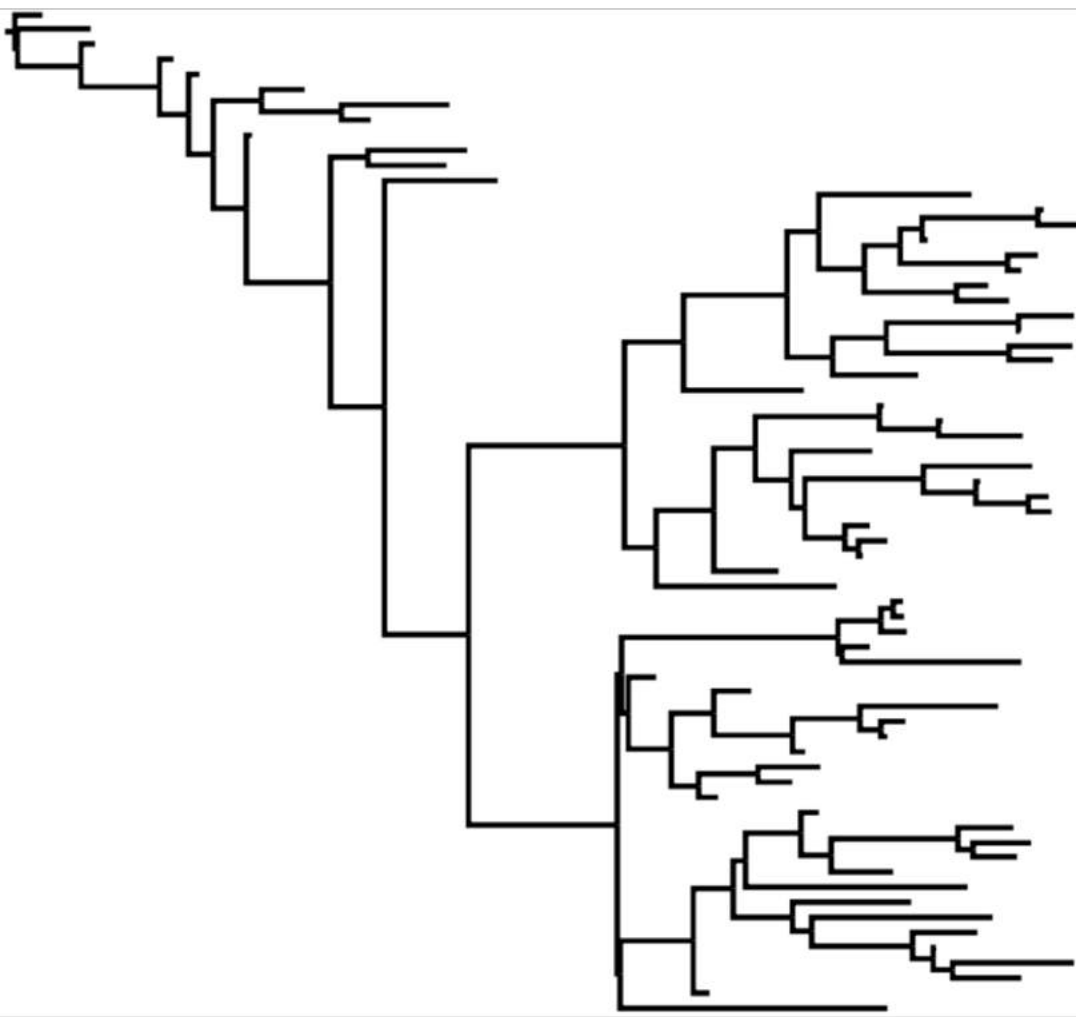
Duchene DA et
al. 2016



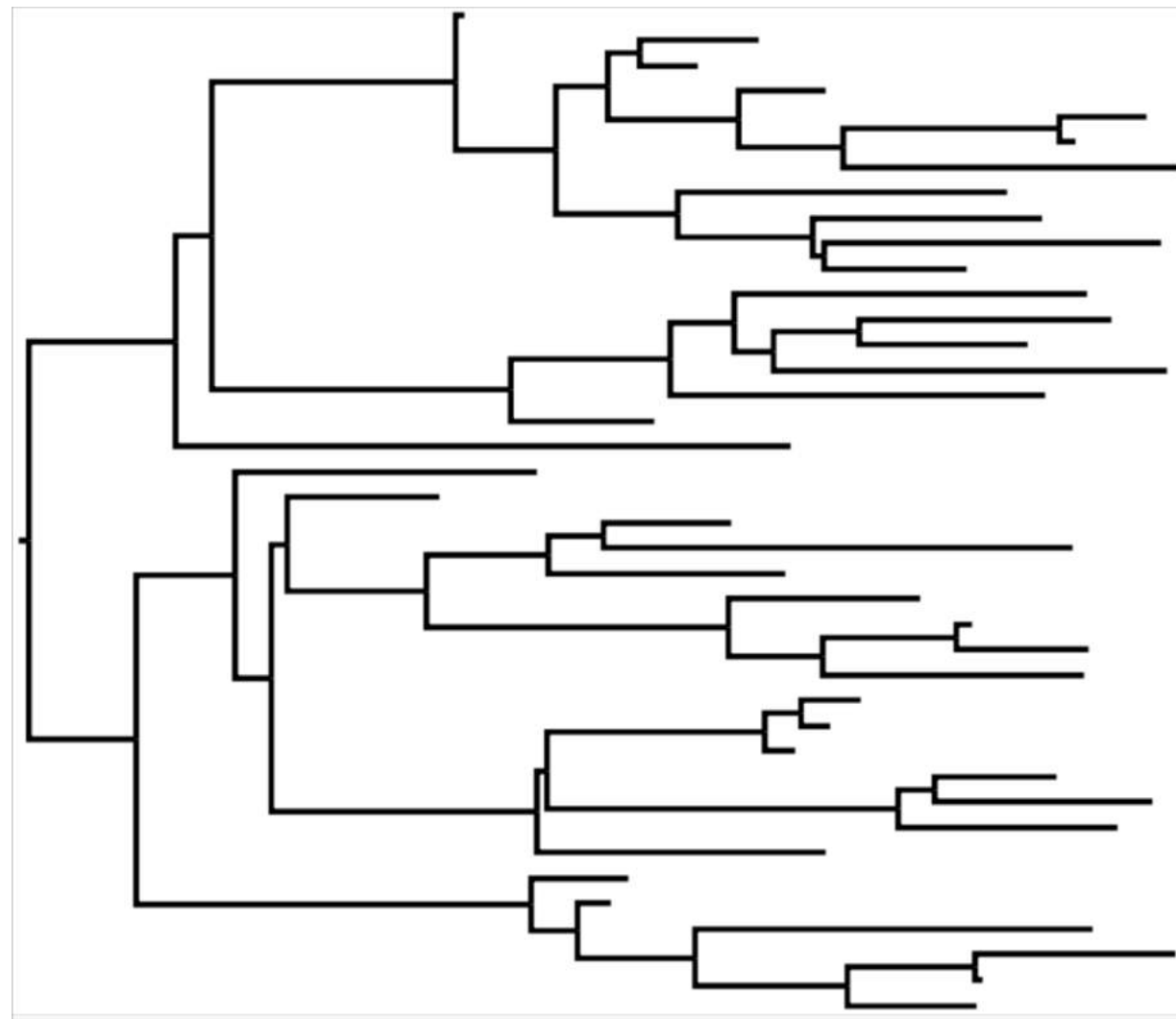
$$P(\text{[Diagram: node with } \lambda_k \text{ and } \delta, p \text{]} \mid \text{[Diagram: tree structure]}) \propto P(\text{[Diagram: tree structure]} \mid \text{[Diagram: node with } \lambda_k \text{ and } \delta, p \text{]}) P(\text{[Diagram: node with } \lambda_k \text{ and } \delta, p \text{]})$$



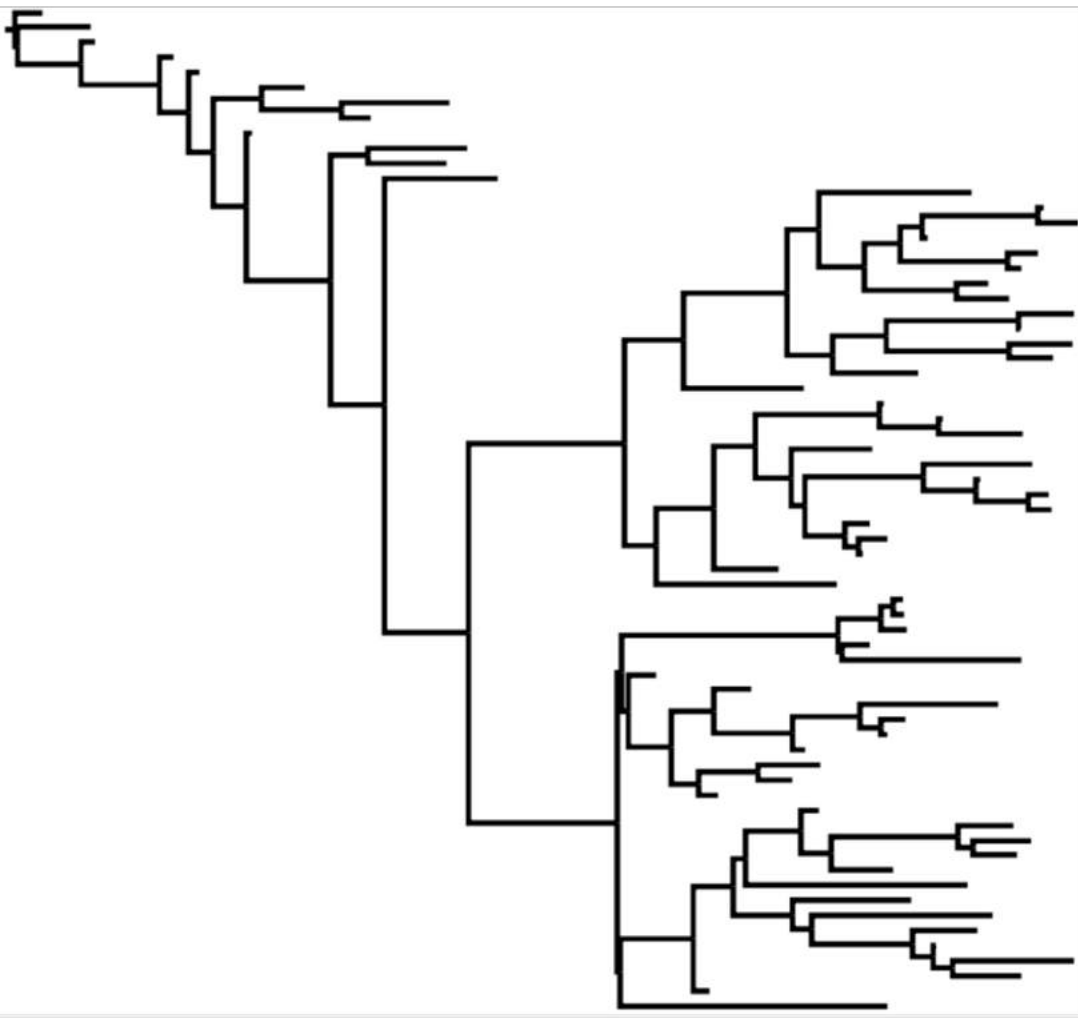
Birth-death constant sampling



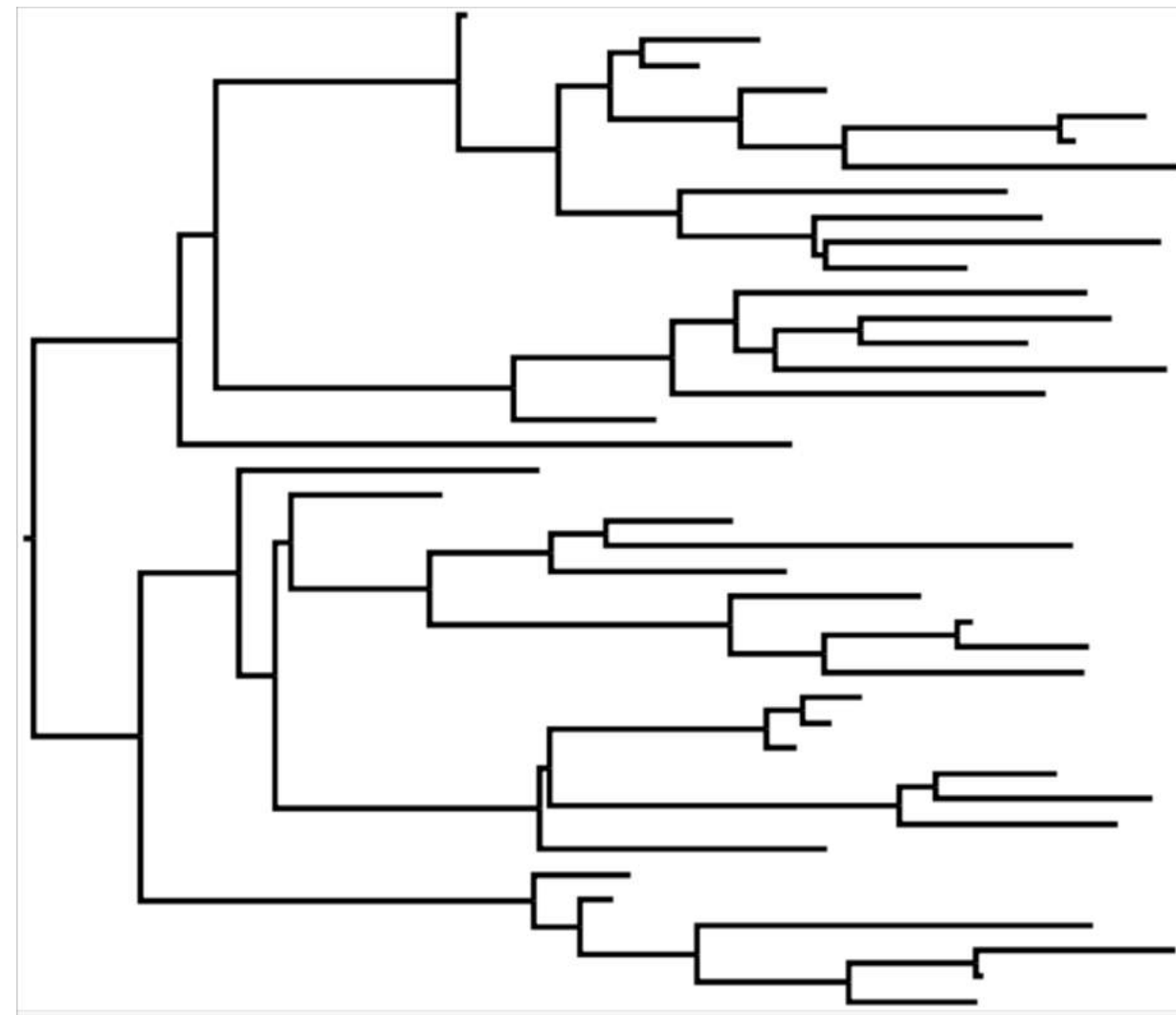
Birth-death skyline- variable sampling



Birth-death constant sampling

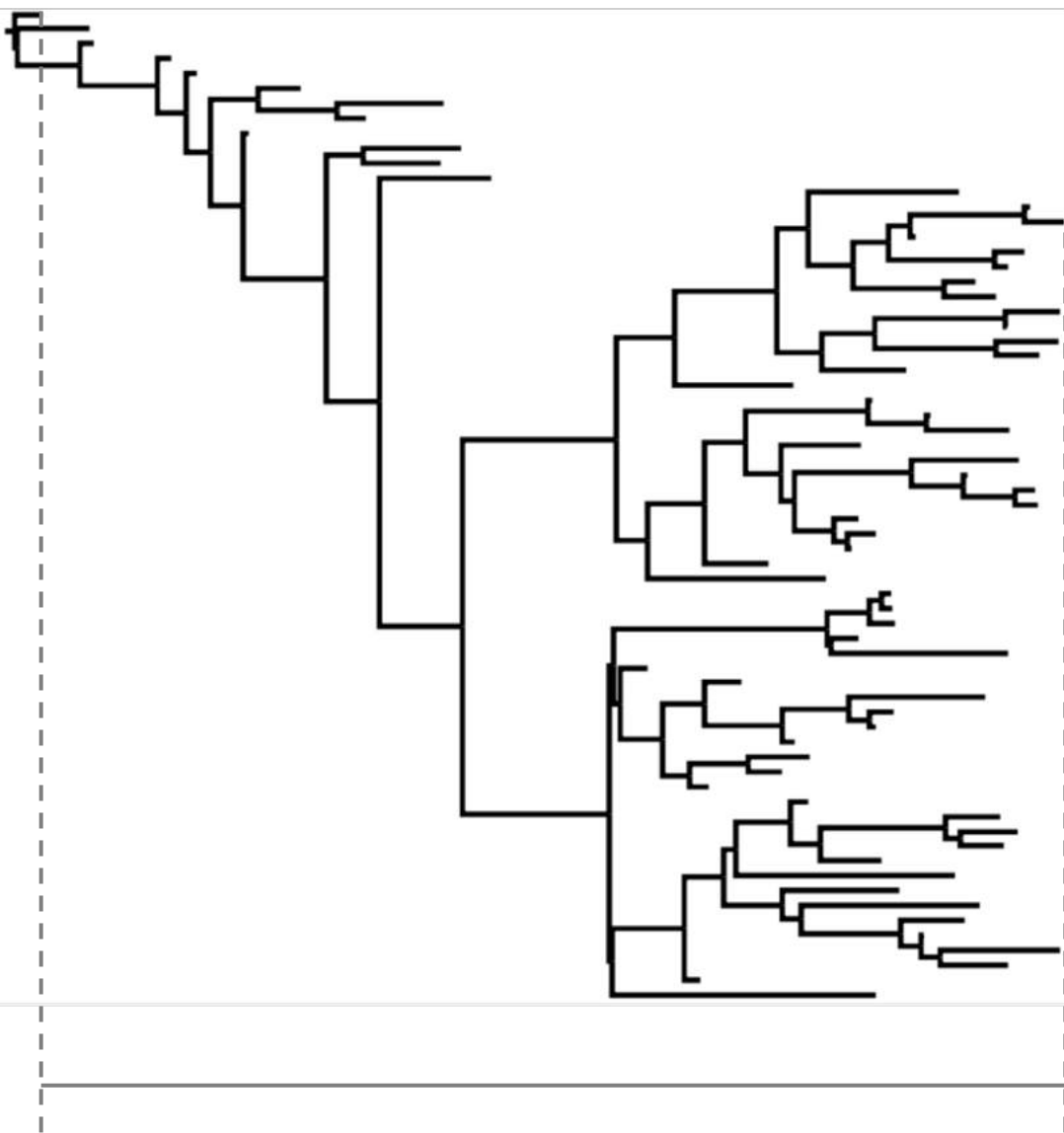


Birth-death skyline- variable sampling



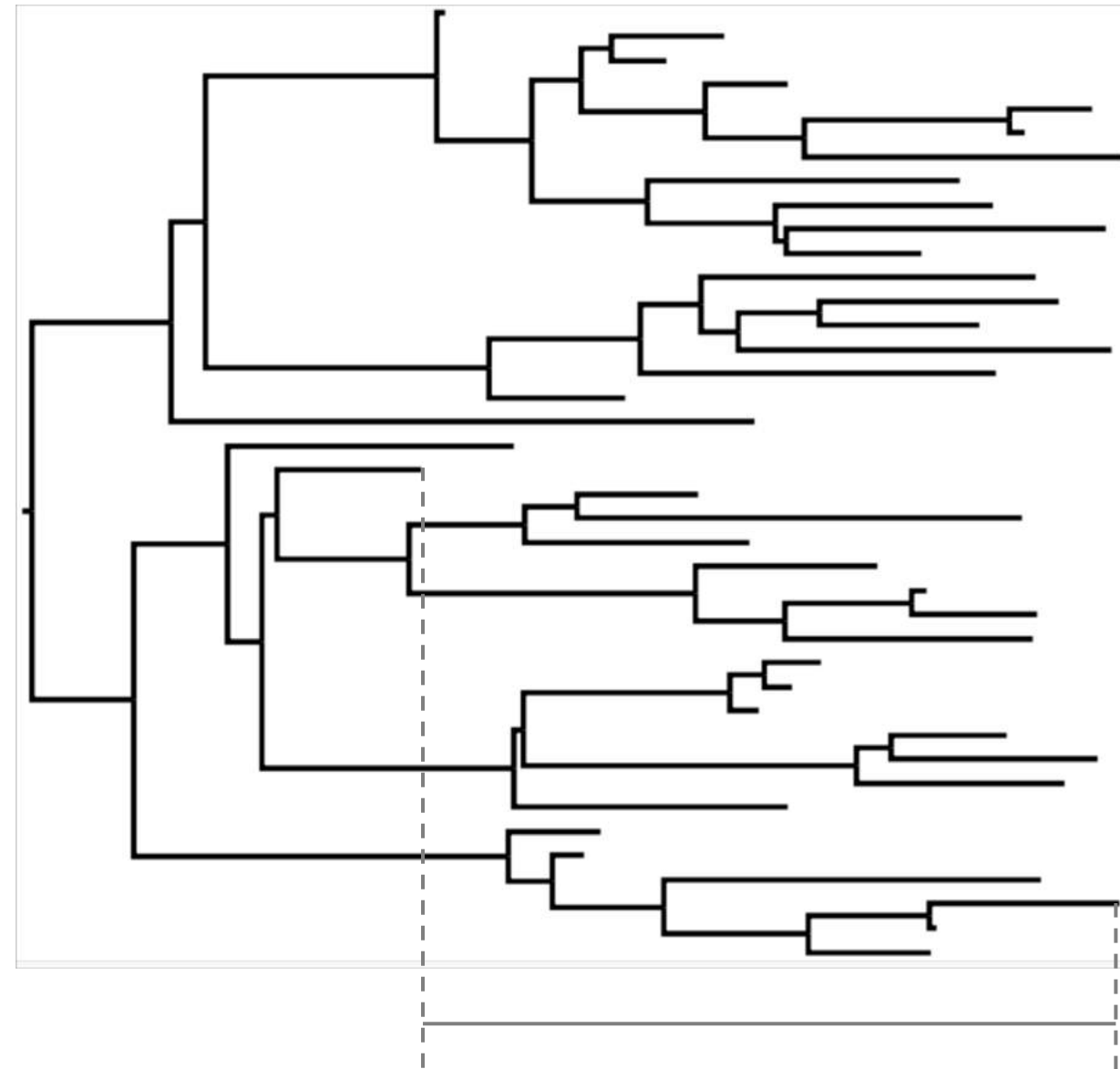
Number of tips

Birth-death constant sampling



10 Months

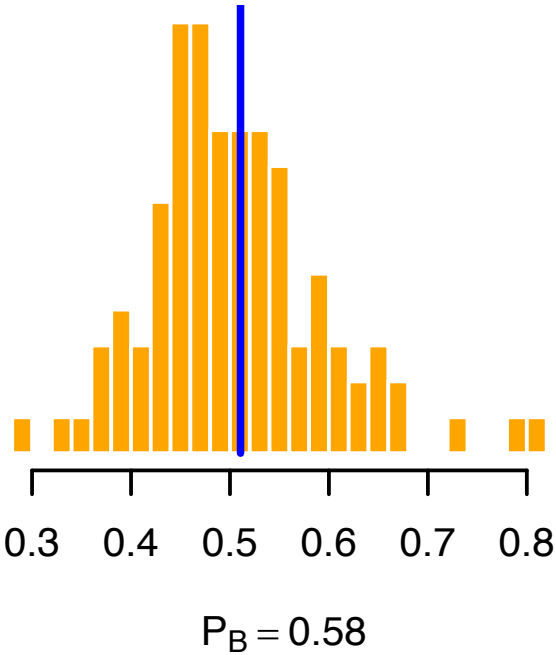
Birth-death skyline- variable sampling



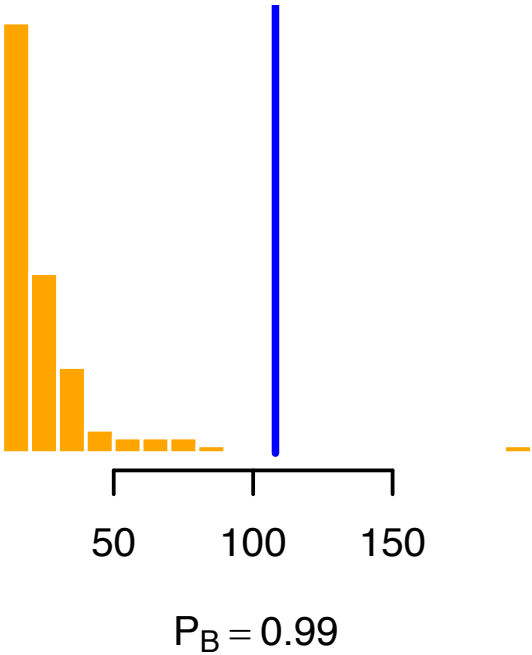
8 Months

Age range of tips

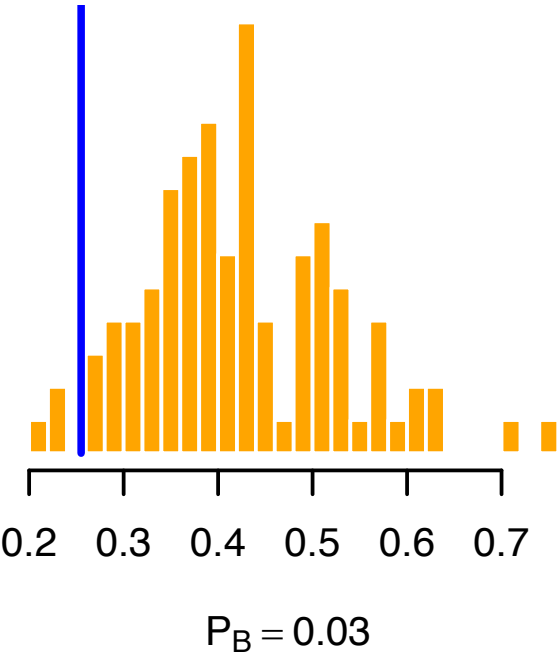
Root height

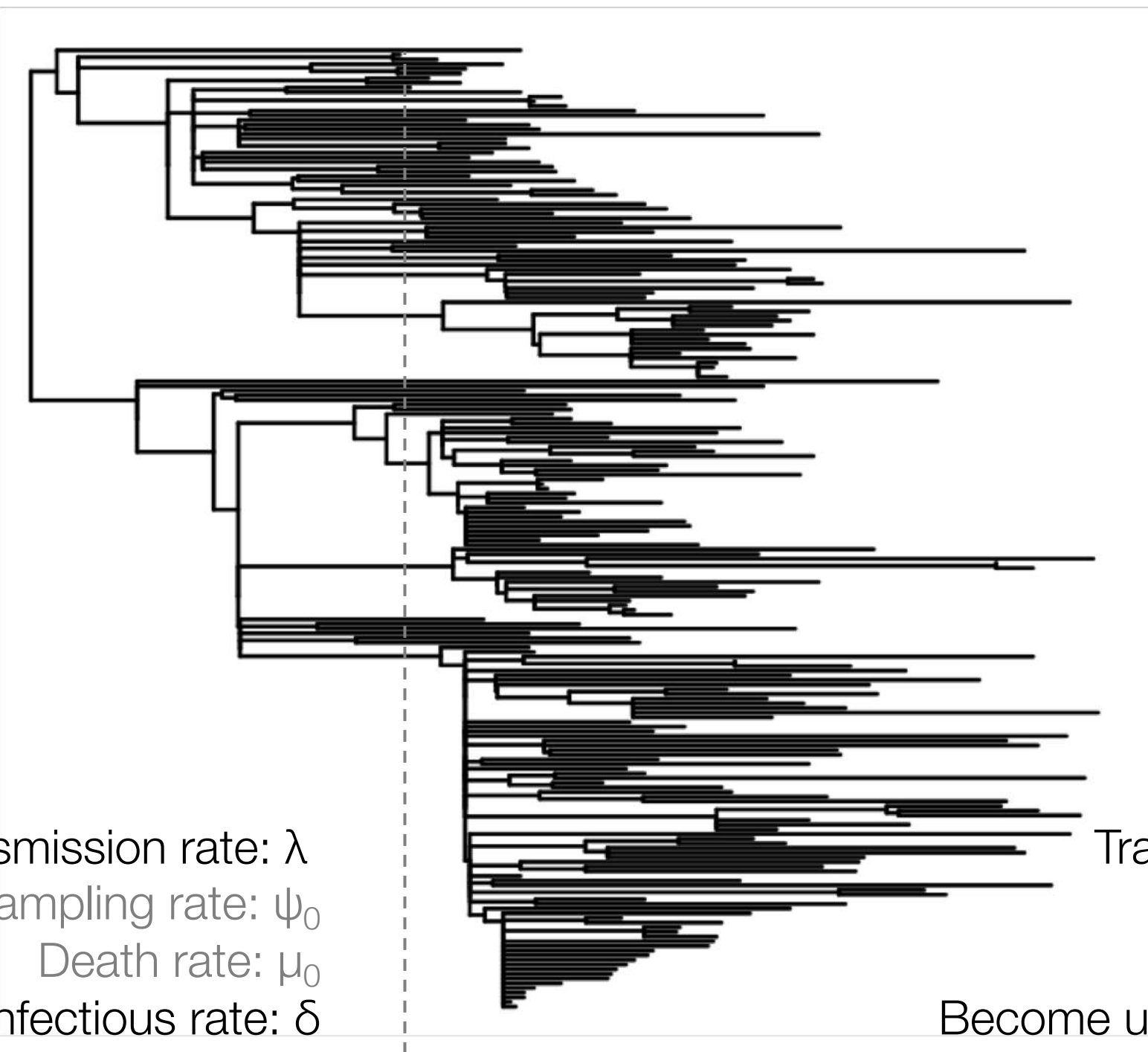


Number of tips



Range of ages
of tips





Transmission rate: λ

Sampling rate: ψ_0

Death rate: μ_0

Become uninfected rate: δ

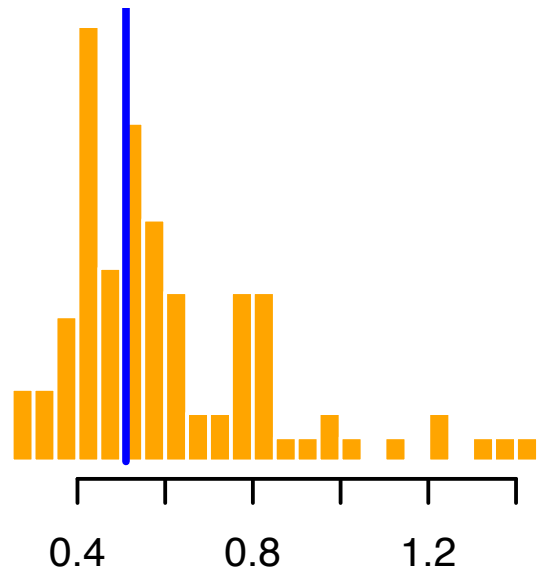
Transmission rate: λ

Sampling rate: ψ_1

Death rate: μ_1

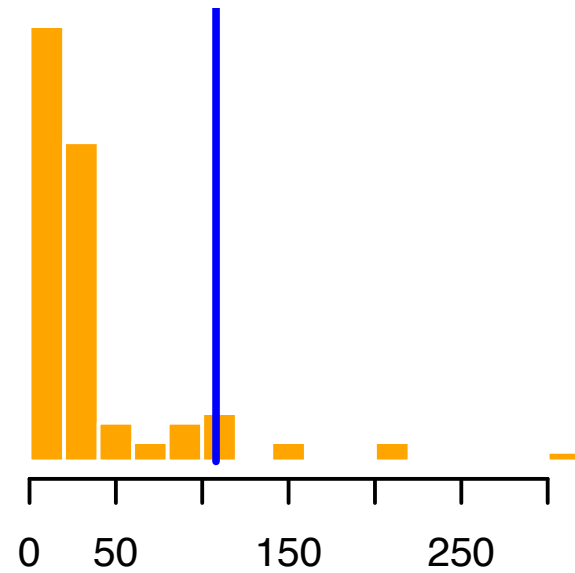
Become uninfected rate: δ

Root height



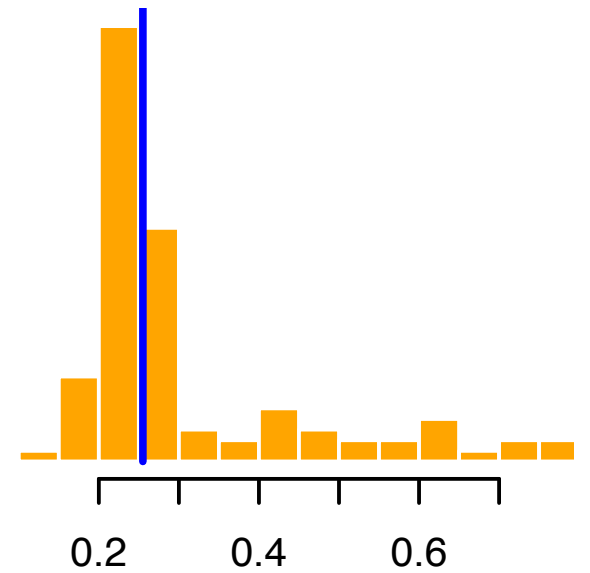
$P_B = 0.42$

Number of tips



$P_B = 0.91$

Range of ages
of tips



$P_B = 0.61$

Final considerations

- Test statistics can be designed using simulations and should assess some expectation or assumption of the model.
- Ideally, test statistics should be a function of the data (certainly not sufficient statistics).
- Comparing the prior and posterior is key in Bayesian model adequacy.
- Model adequacy can be useful for model improvement and for understanding the reliability of the inferences.