

0111010  
0101010  
0001010  
0111011  
0111010  
0111010

# Infectious disease phylodynamics

Louis du Plessis

**ETH** zürich

**DBSSE**

# On the program for today

---

## **(1) What is phylodynamics?**

- (2) Bayesian inference recap
- (3) BEAST2 introduction

## **Tutorial:** Molecular clock dating (part i)

- (4) Molecular clock models

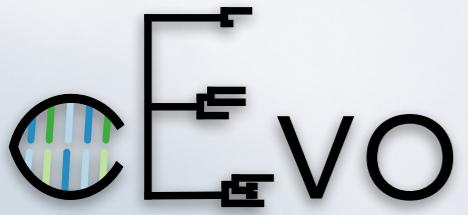
## **Tutorial:** Molecular clock dating (part ii)

- (5) Setting priors

## **Tutorial:** Phylodynamics (part i)

- (6) Tree priors

## **Tutorial:** Phylodynamics (part ii)



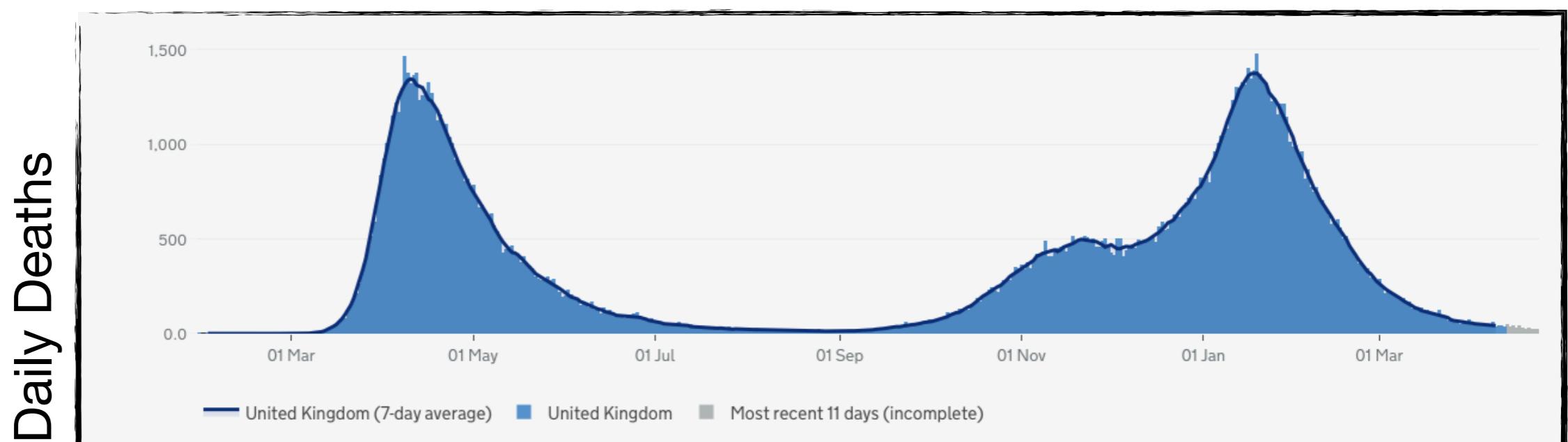
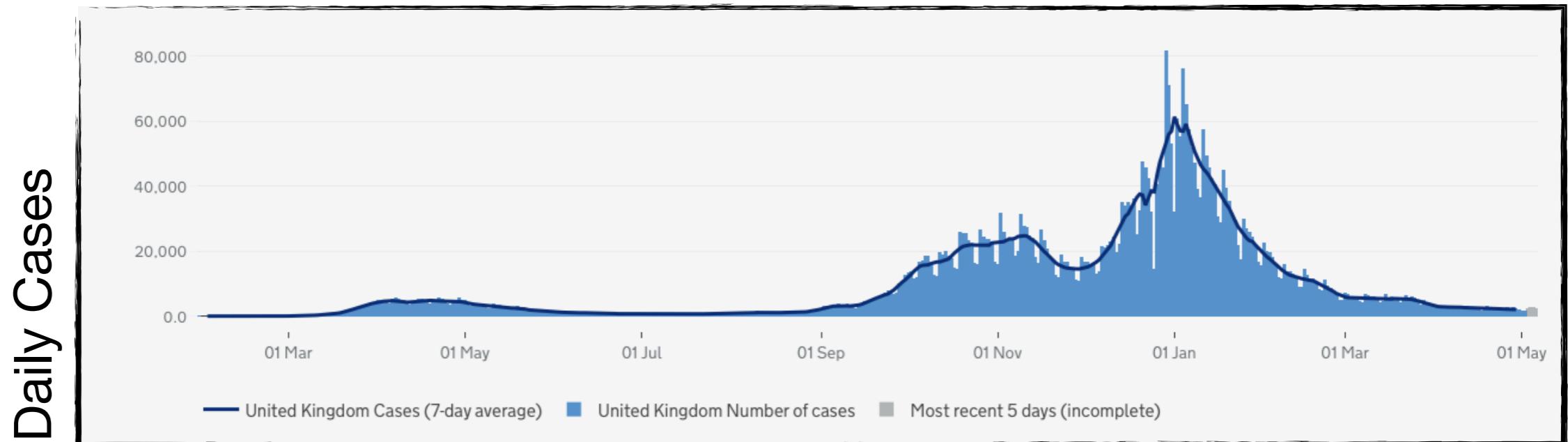
# What is phylodynamics?

## Louis du Plessis

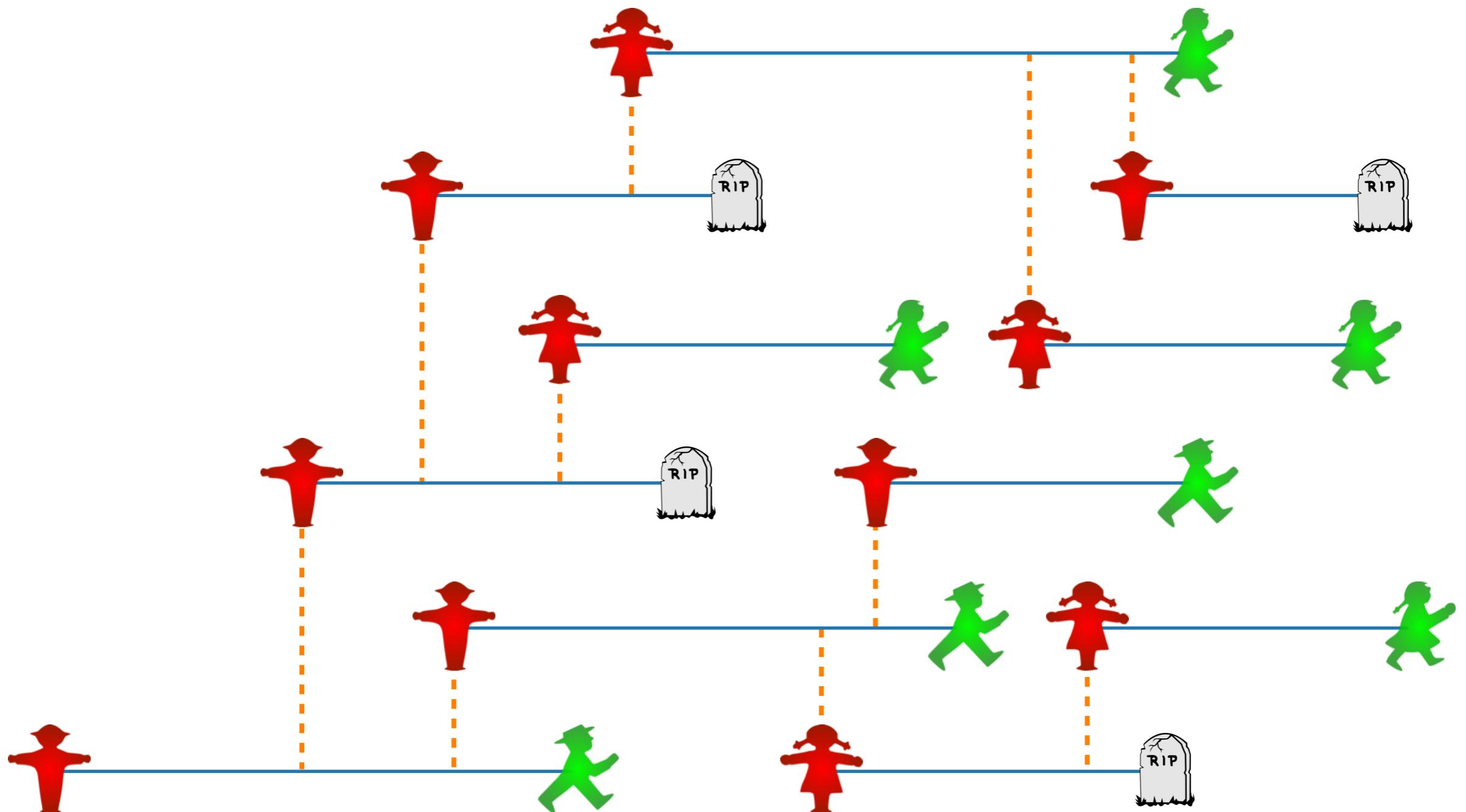
ETH zürich

DBSSE

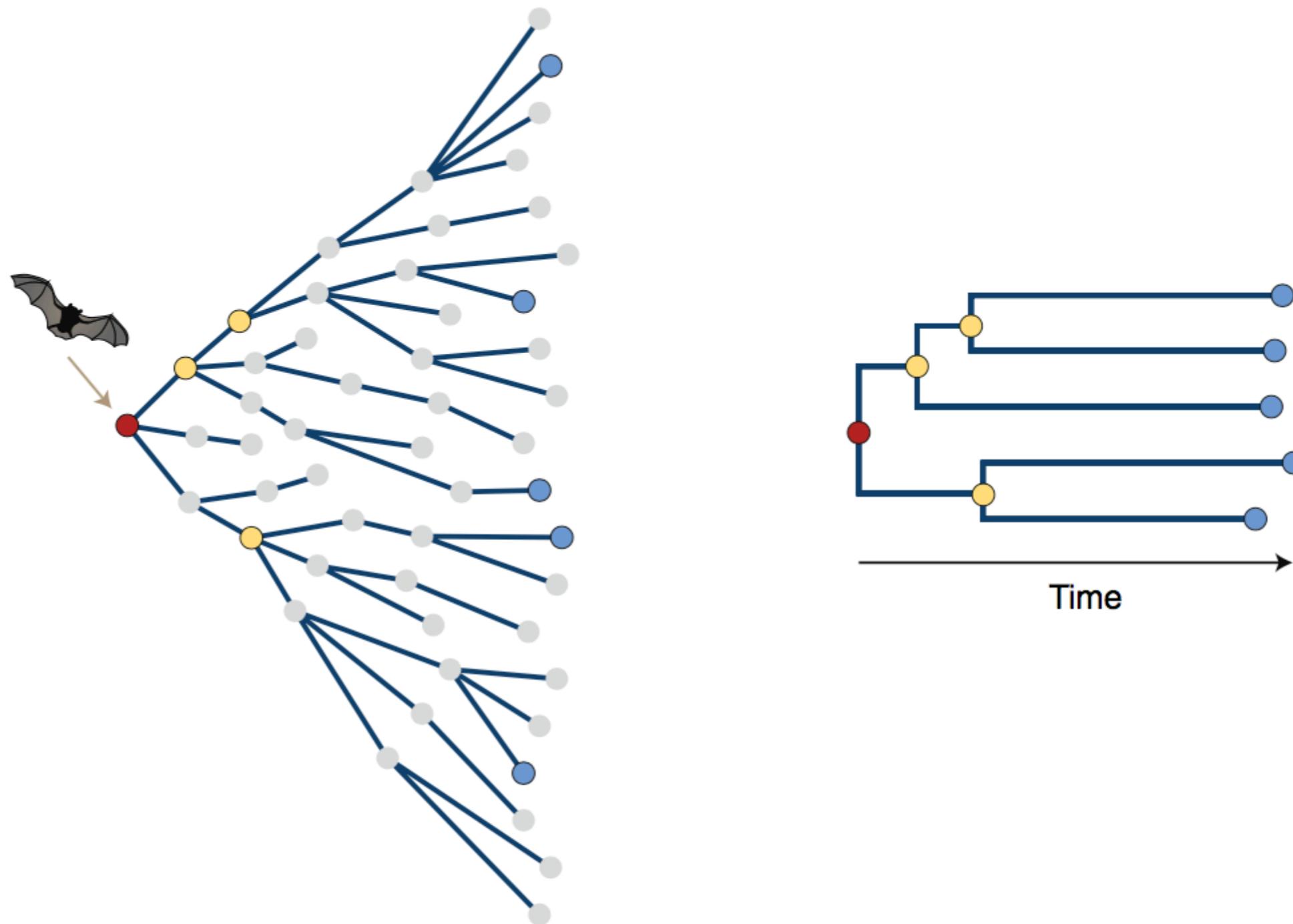
# Epidemic time-series (United Kingdom)

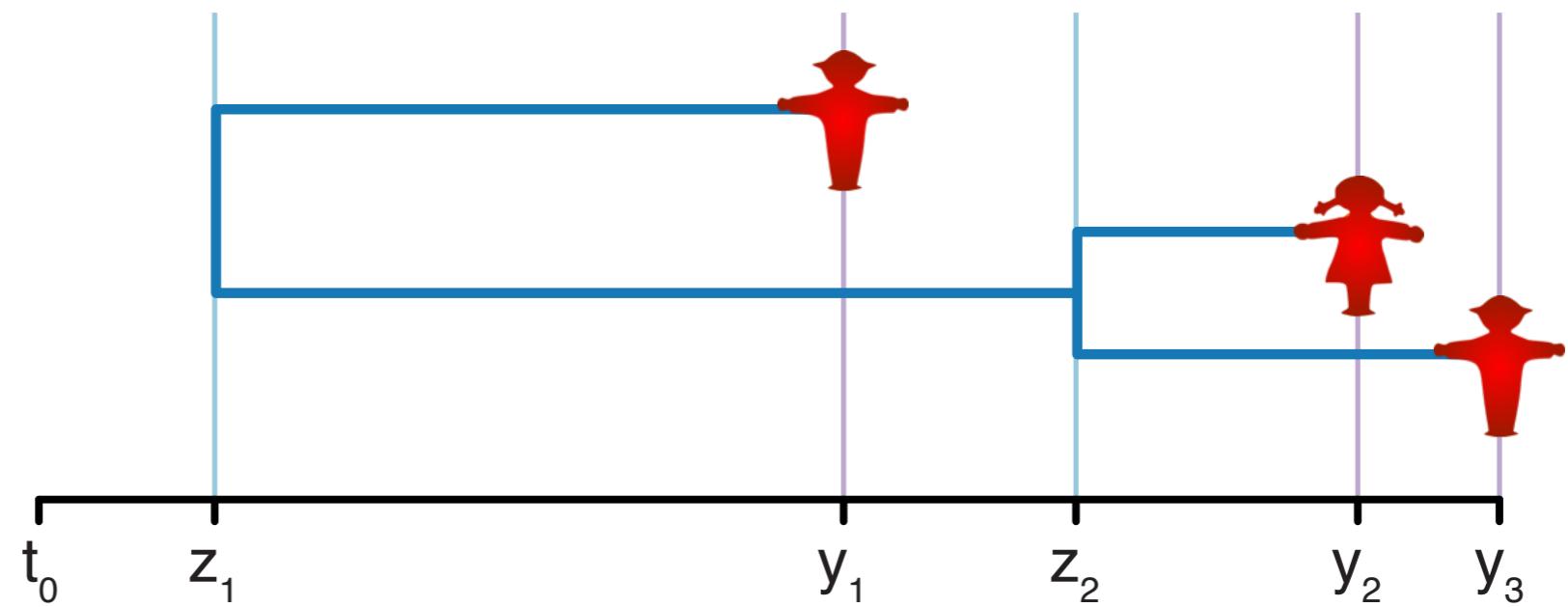
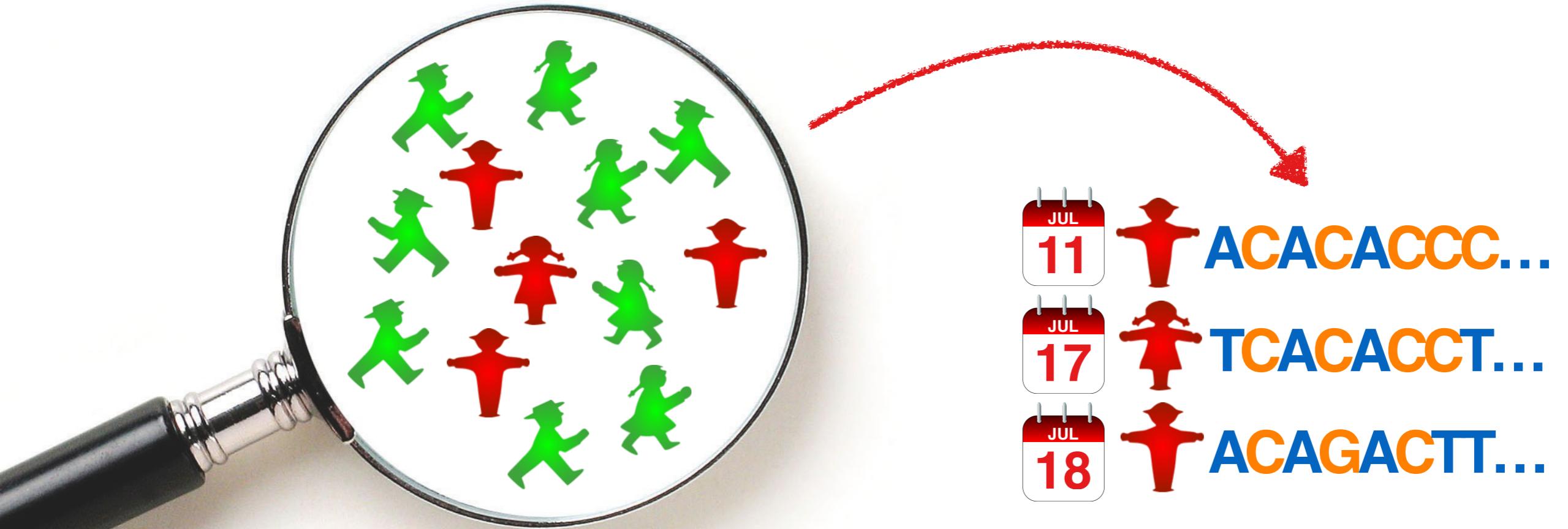


# Definition: Transmission tree



# Definition: Partial transmission tree (phylogeny)

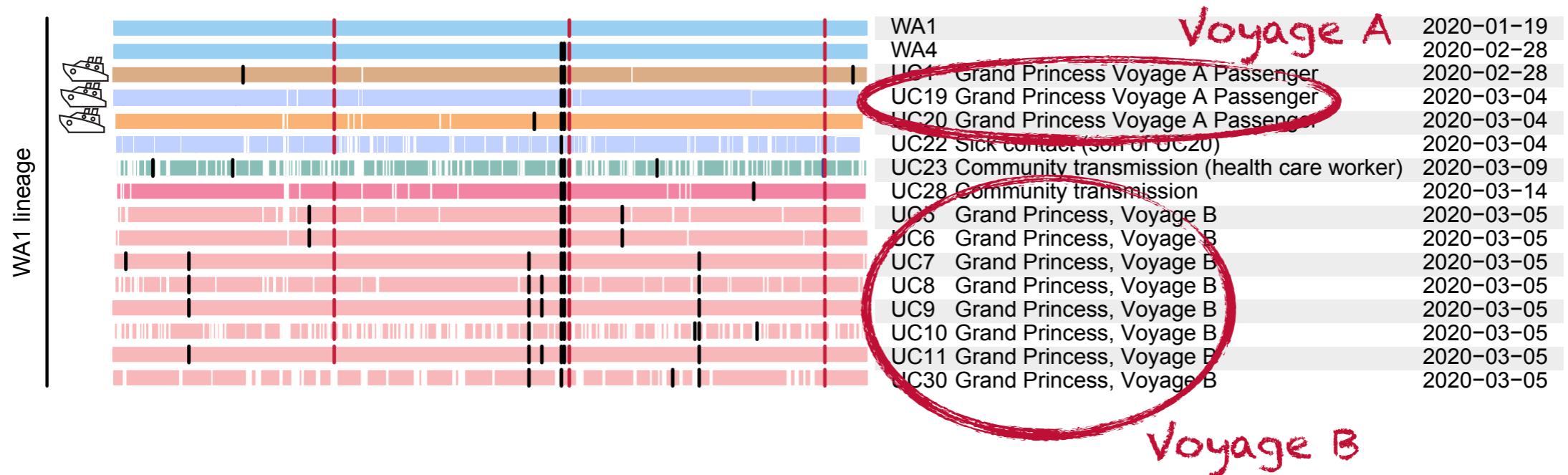
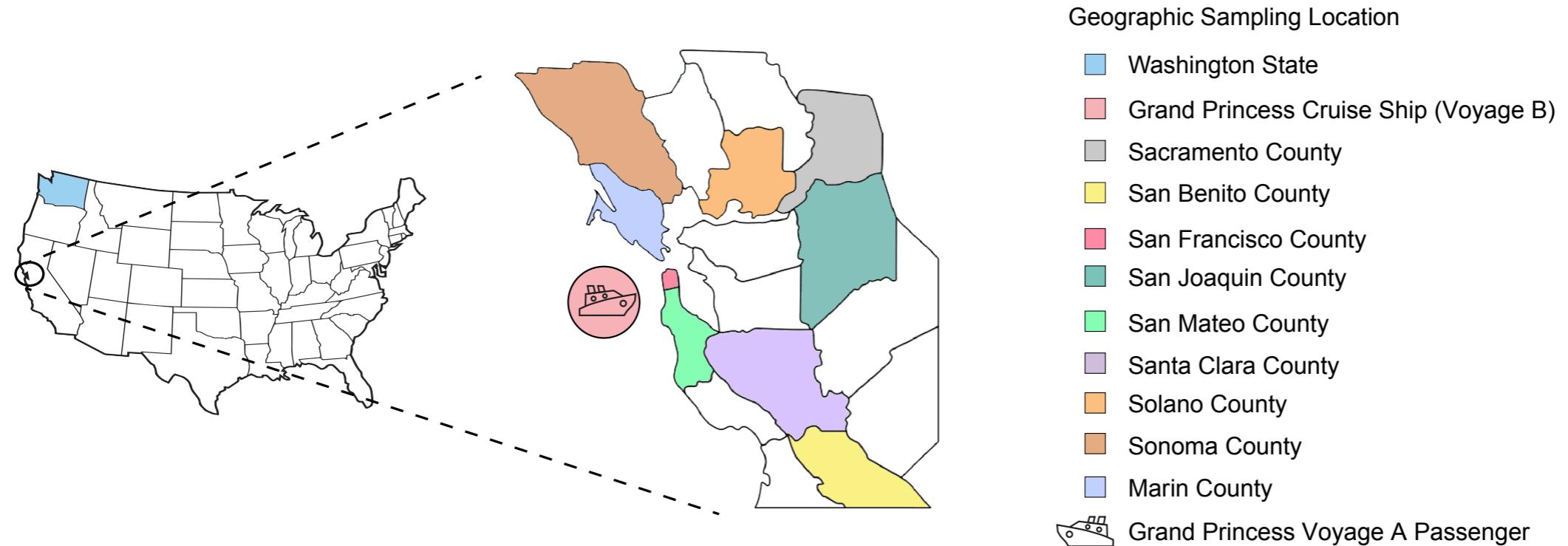






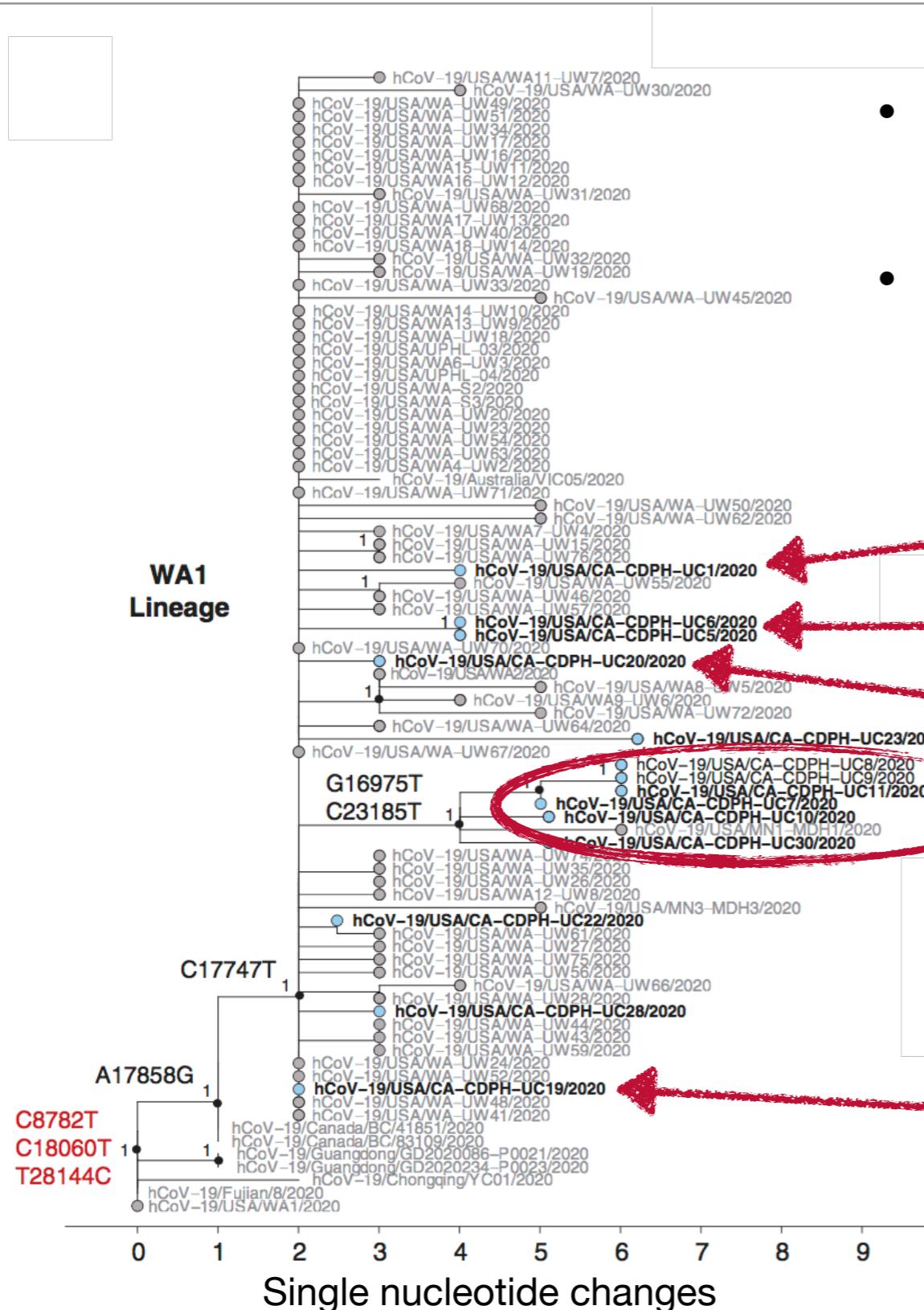
# SARS-CoV-2 in Northern California

(Feb - Mar 2020)



# SARS-CoV-2 in Northern California

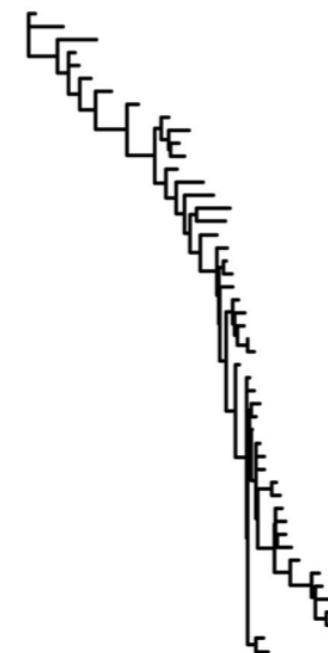
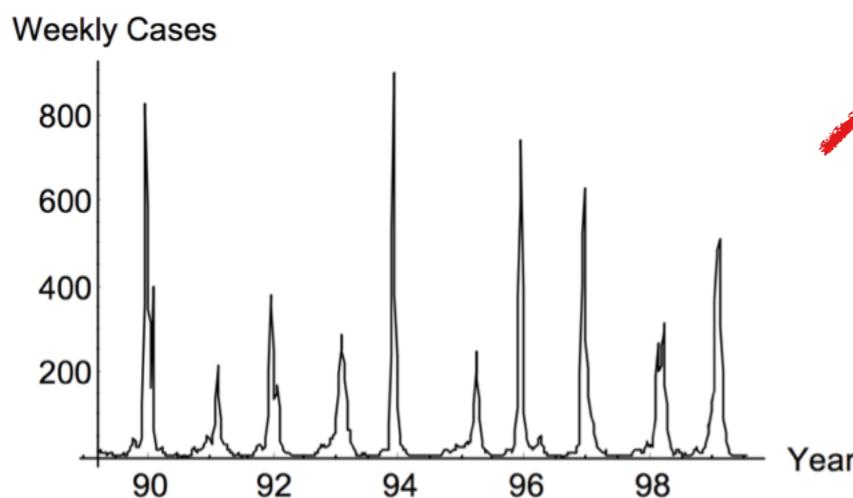
(Feb - Mar 2020)



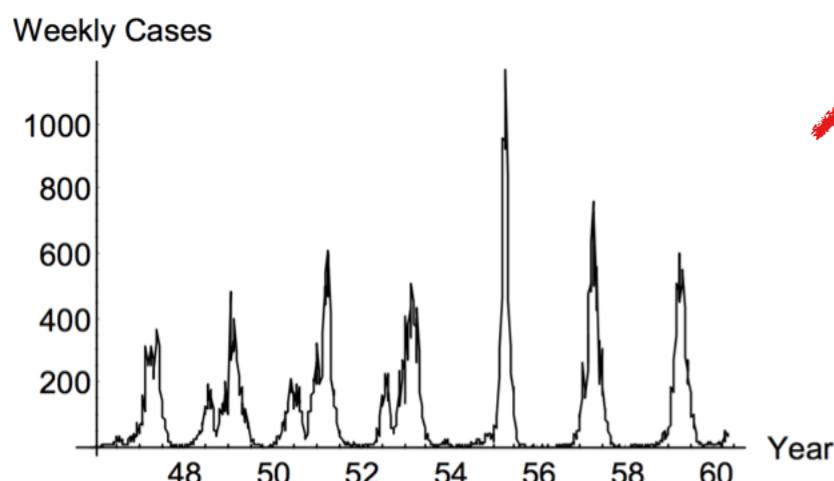
- Transmission from **Voyage A** (11 - 21 Feb) to **Voyage B** (22 Feb - 4 Mar)
- Community transmission of the WA1 lineage in California by **11 Feb**

# Case data alone is not enough!

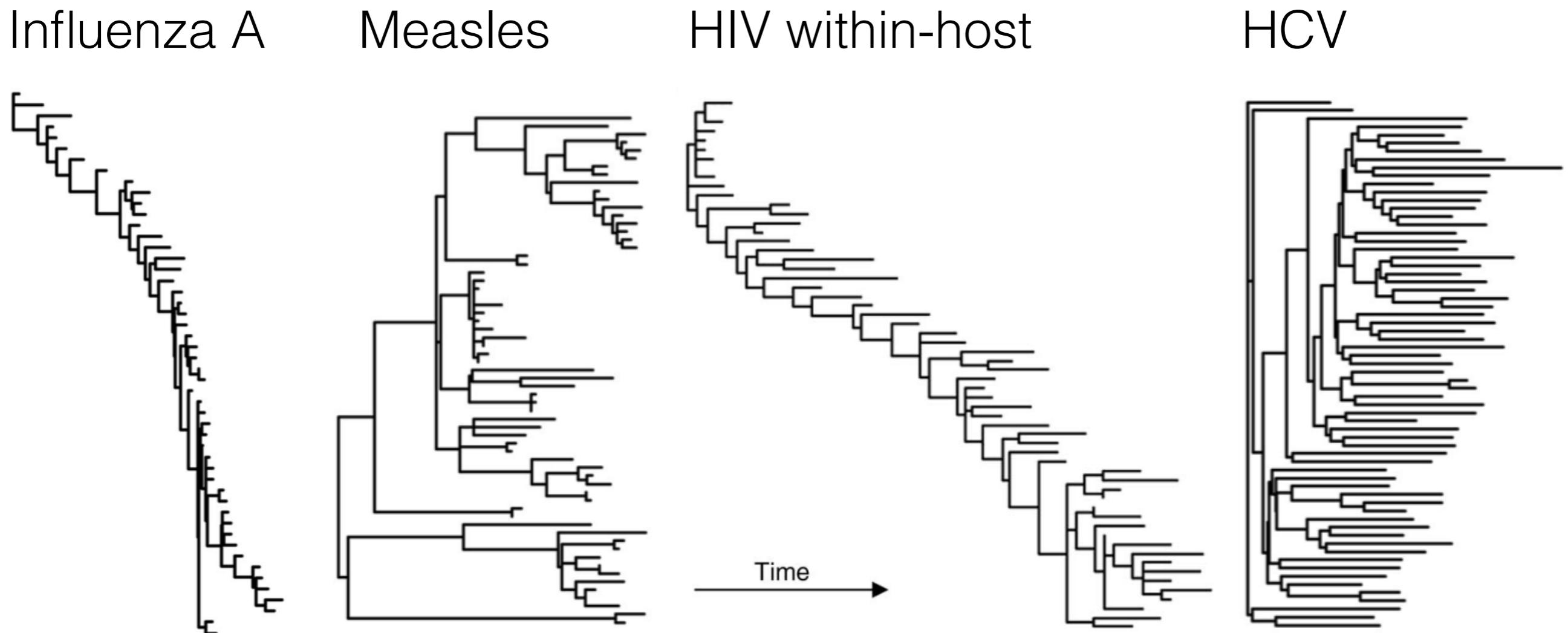
## Influenza A



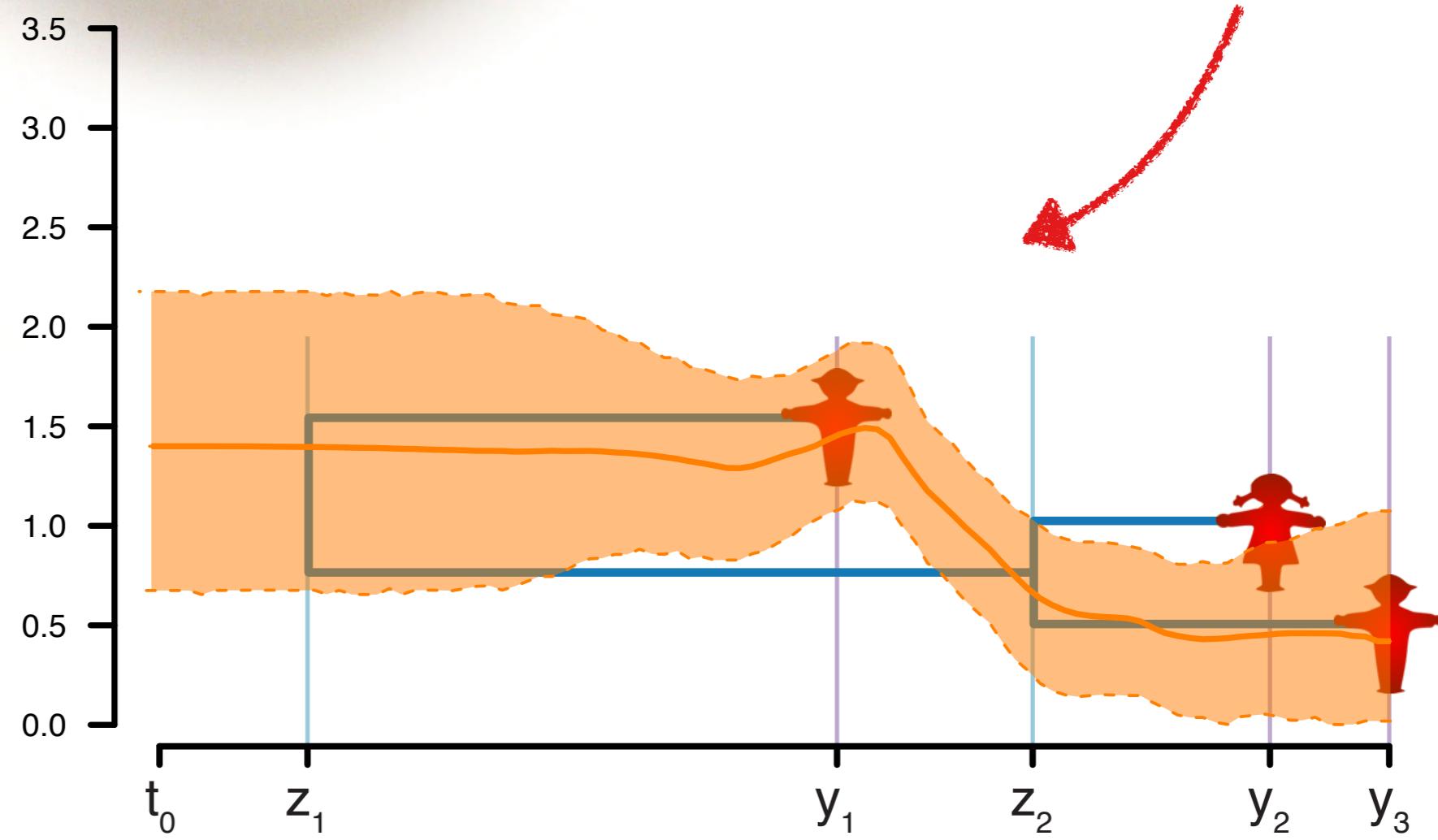
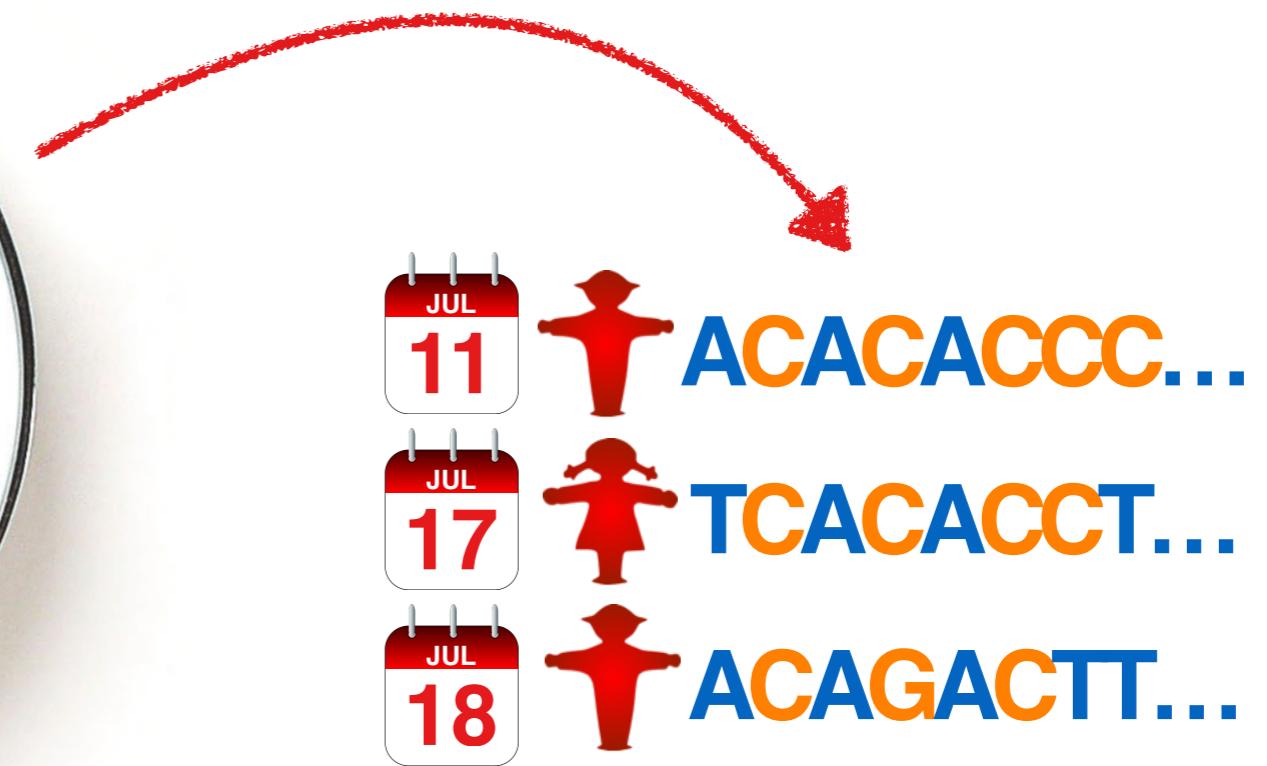
## Measles



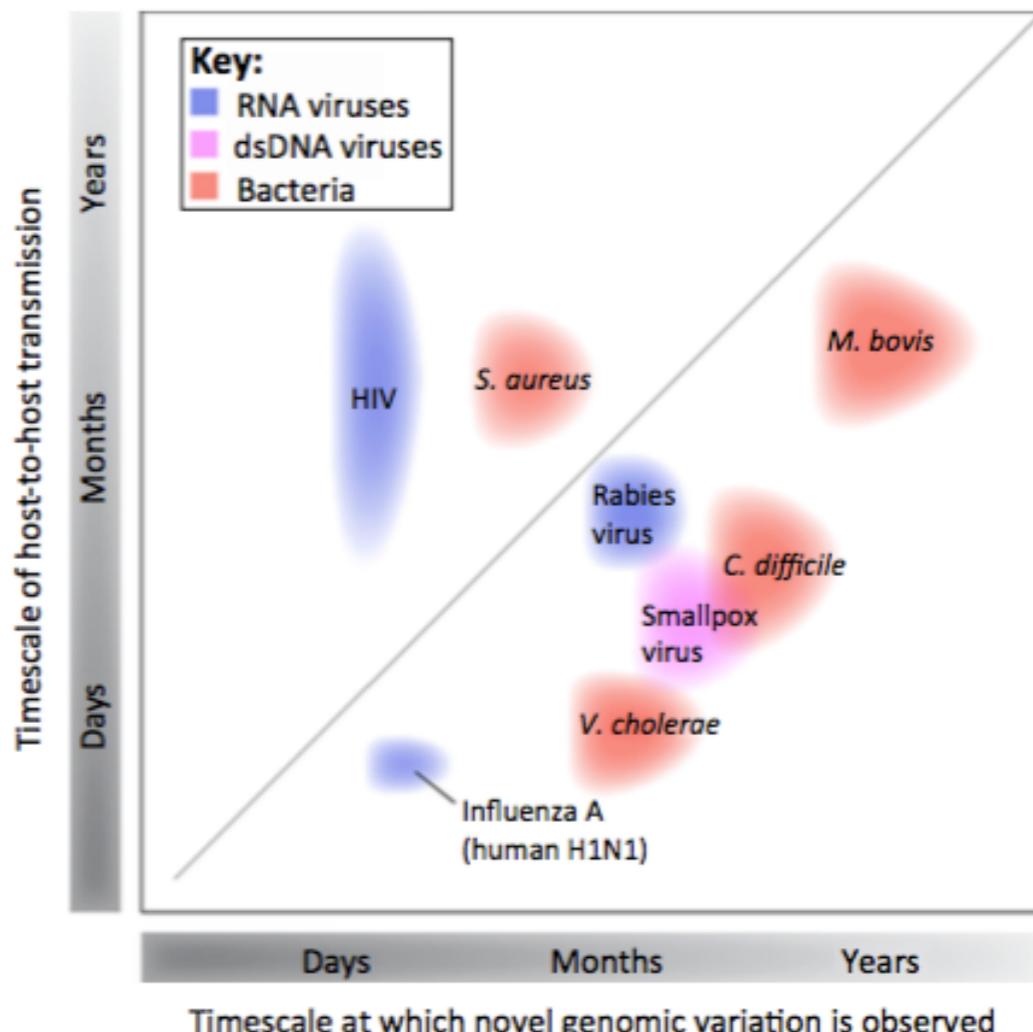
# Genomes contain a signature of the epidemiological dynamics



**Phyldynamics** aims to infer the **dynamics** responsible for the observed **phylogenies**



# Measurably evolving populations

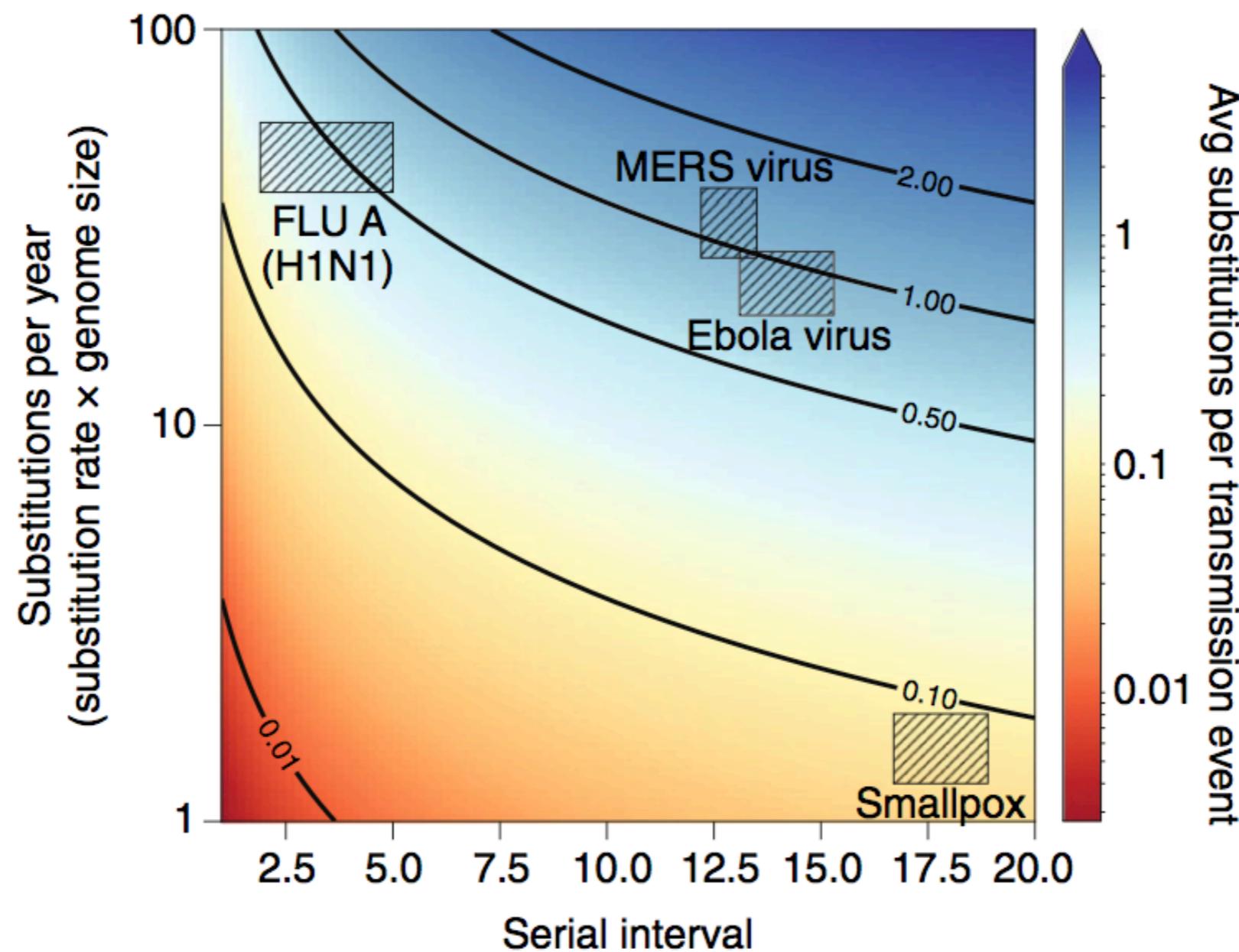


**Novel genetic variation** accumulates over the **sampling period**

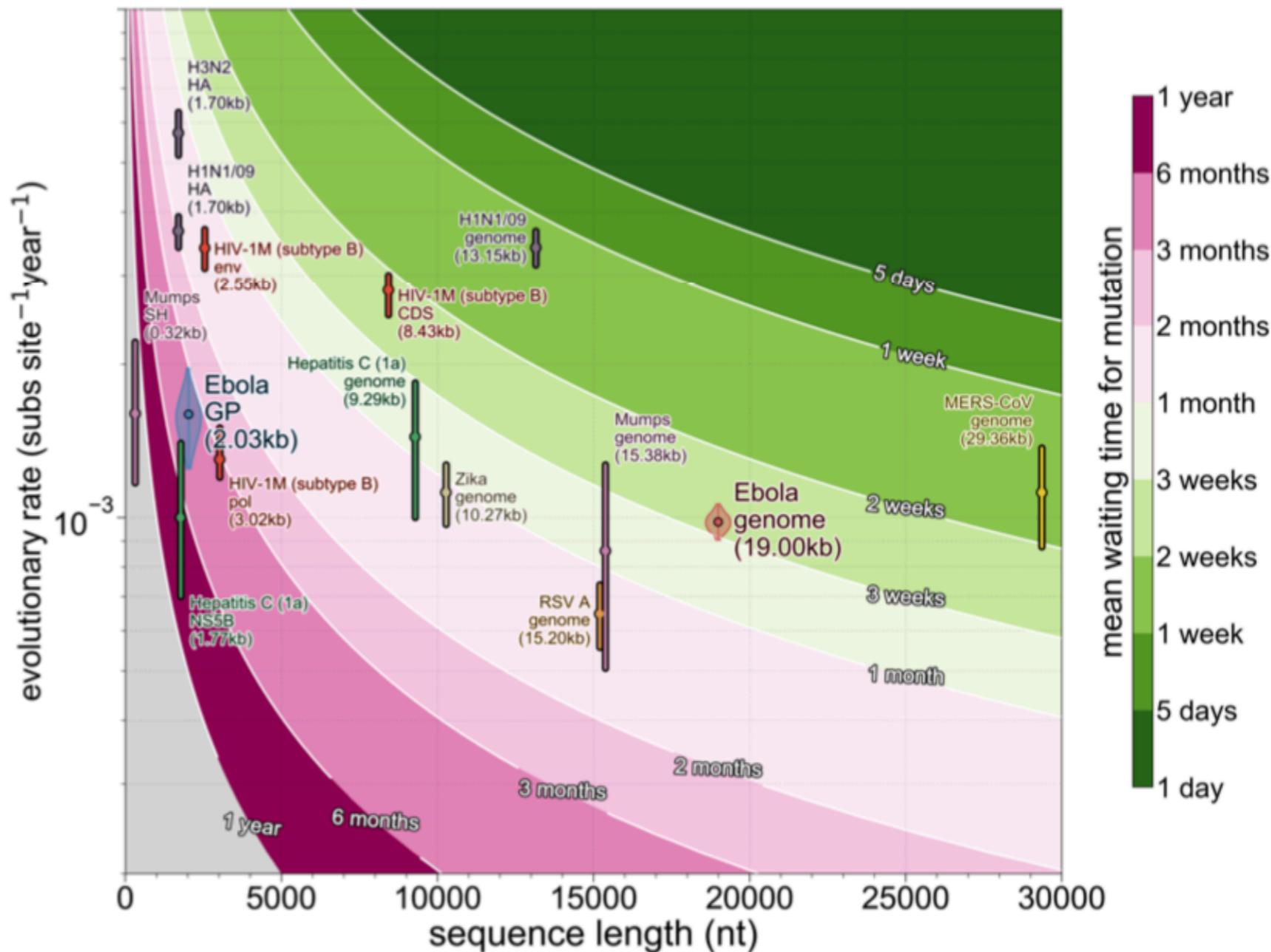
- Large population size
- High mutation rate
- Short generation times
- This is the case for many pathogens (especially viruses)

**Epidemiological** and **evolutionary** dynamics occur on the same **timescale!**

# Measurably evolving populations and molecular clocks



# Measurably evolving populations and molecular clocks

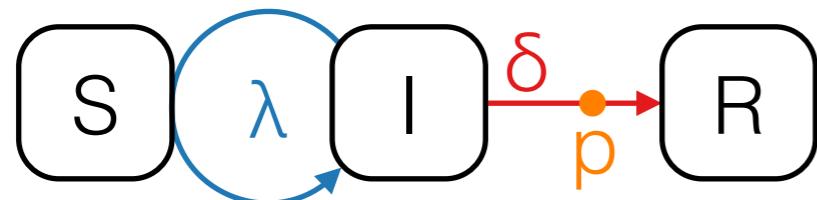
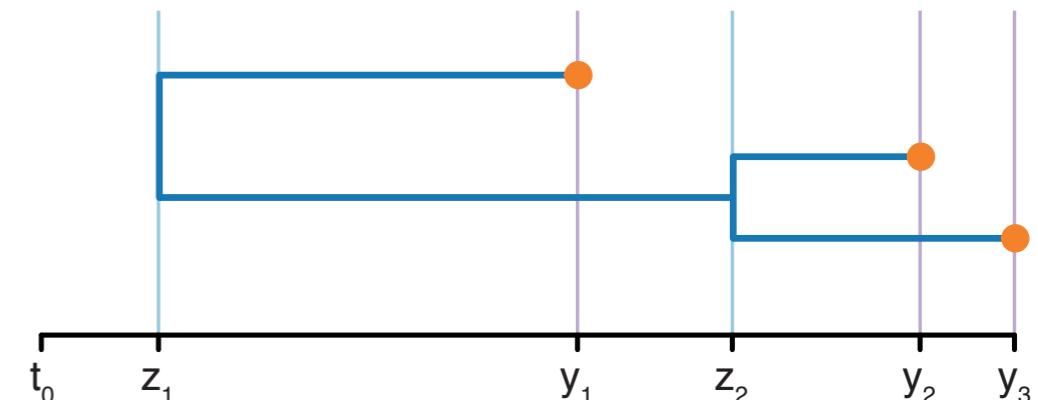


# What is phylodynamics?

---

## Phylogenetics

- State of process
- Classification
- What are the ancestral relationships?



## Phylodynamics

- Dynamics of process
- Epidemiological parameters
- How did we get here?

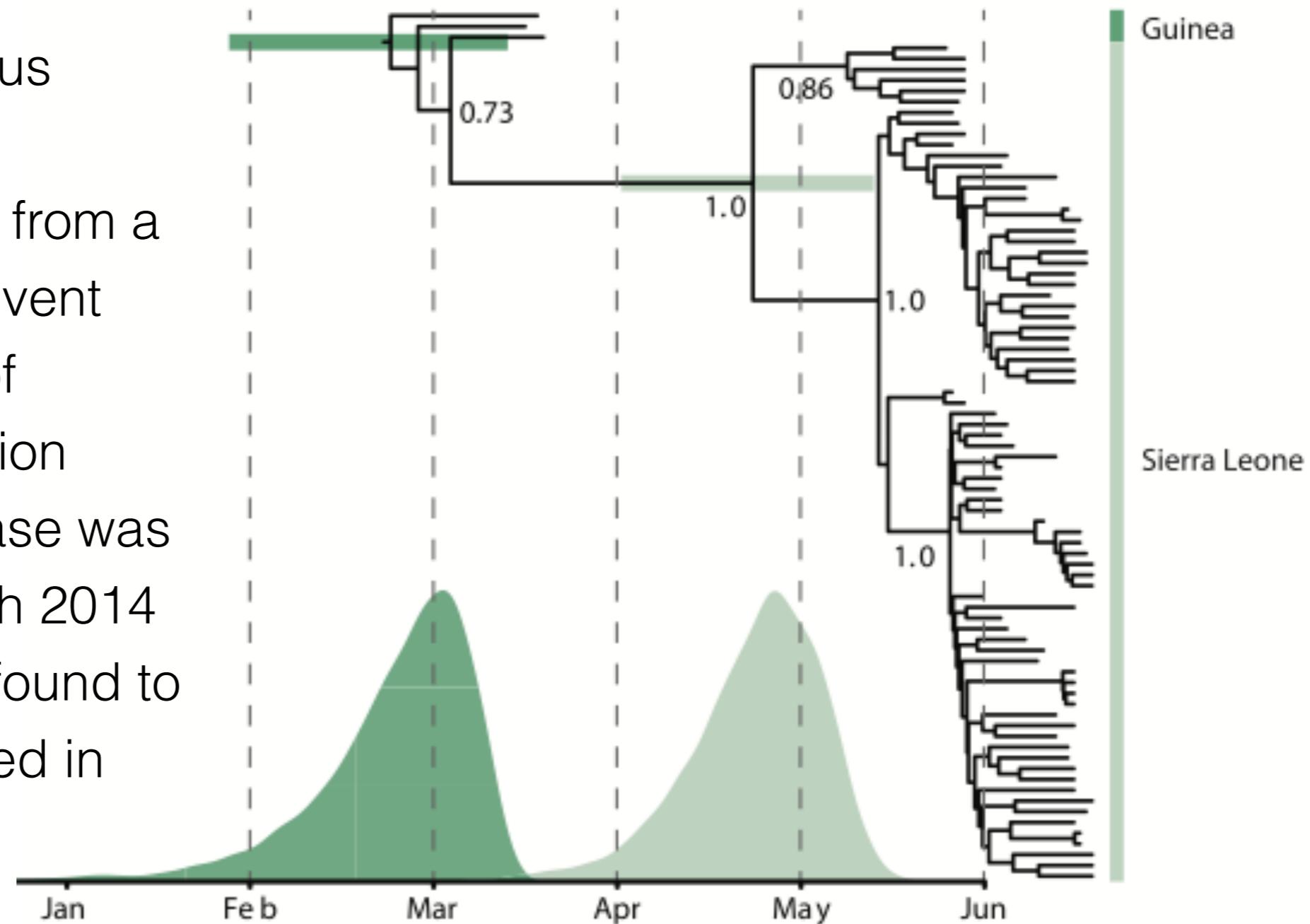
# Phyldynamics questions

---

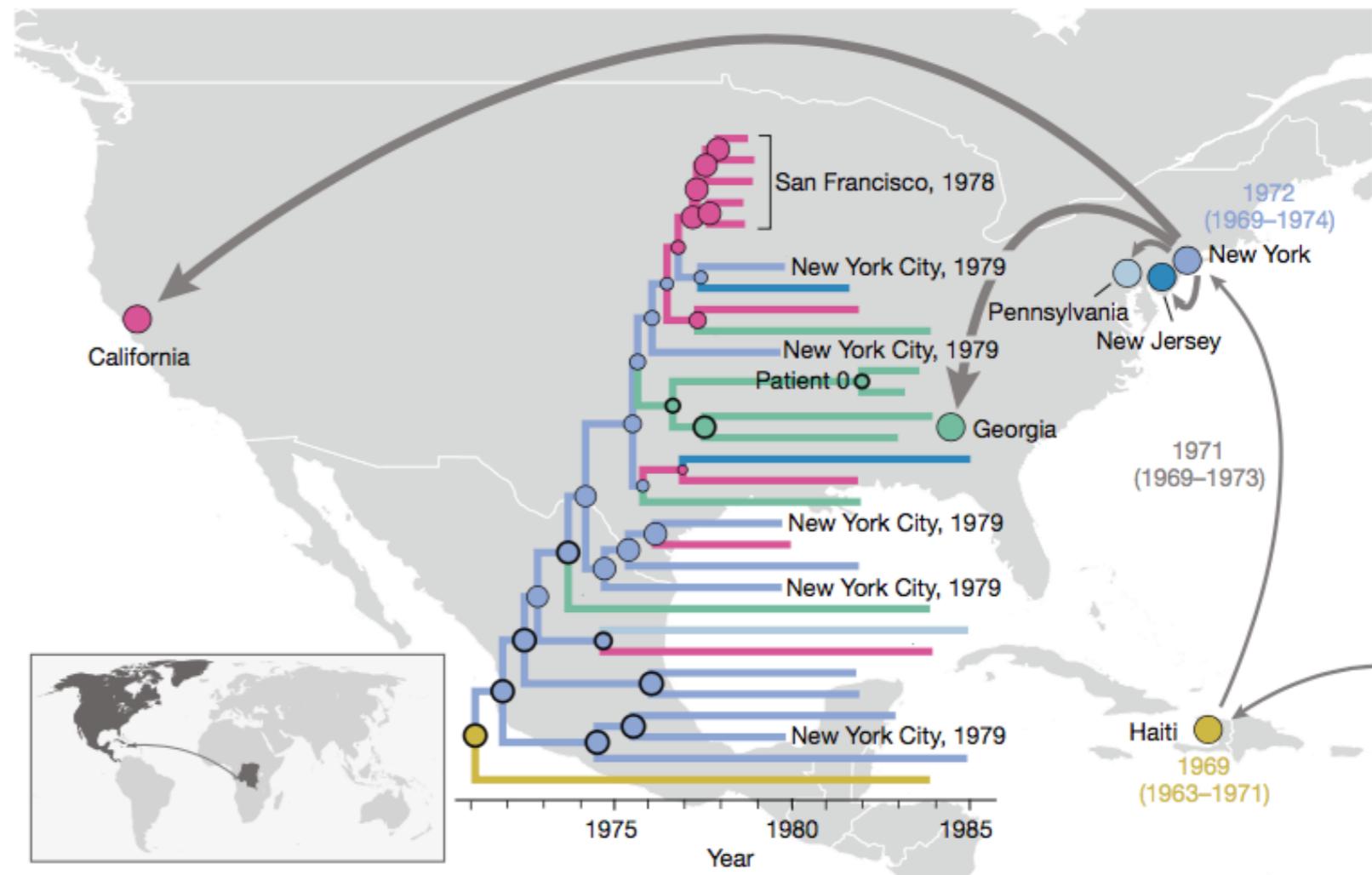
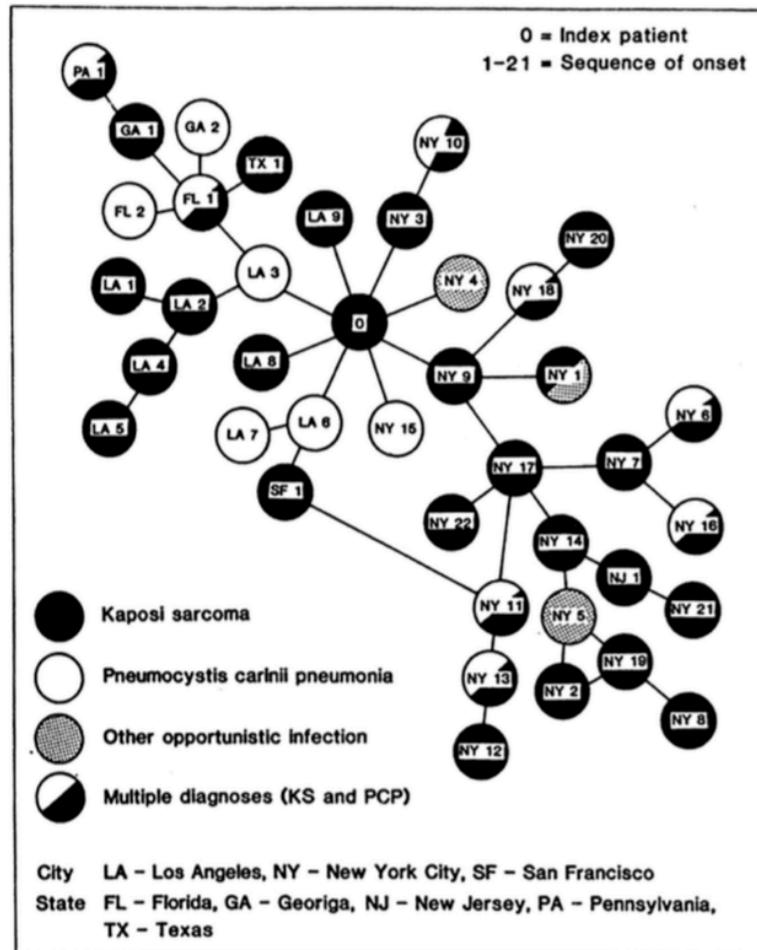
- How genetically diverse is a pathogen population?
- How do pathogen genomes change through time?
- How does pathogen genetic diversity vary through space?
- What processes or events determine these changes?
- When did an epidemic start?
- Where did it come from?
- How fast is it transmitting?
- In what direction is it spreading?
- Are hosts X,Y & Z epidemiologically linked?

# West African Ebola virus epidemic

- Biggest Ebola virus epidemic
- Epidemic started from a single zoonosis event
- Several months of cryptic transmission before the first case was detected in March 2014
- Index case later found to have been infected in December 2013



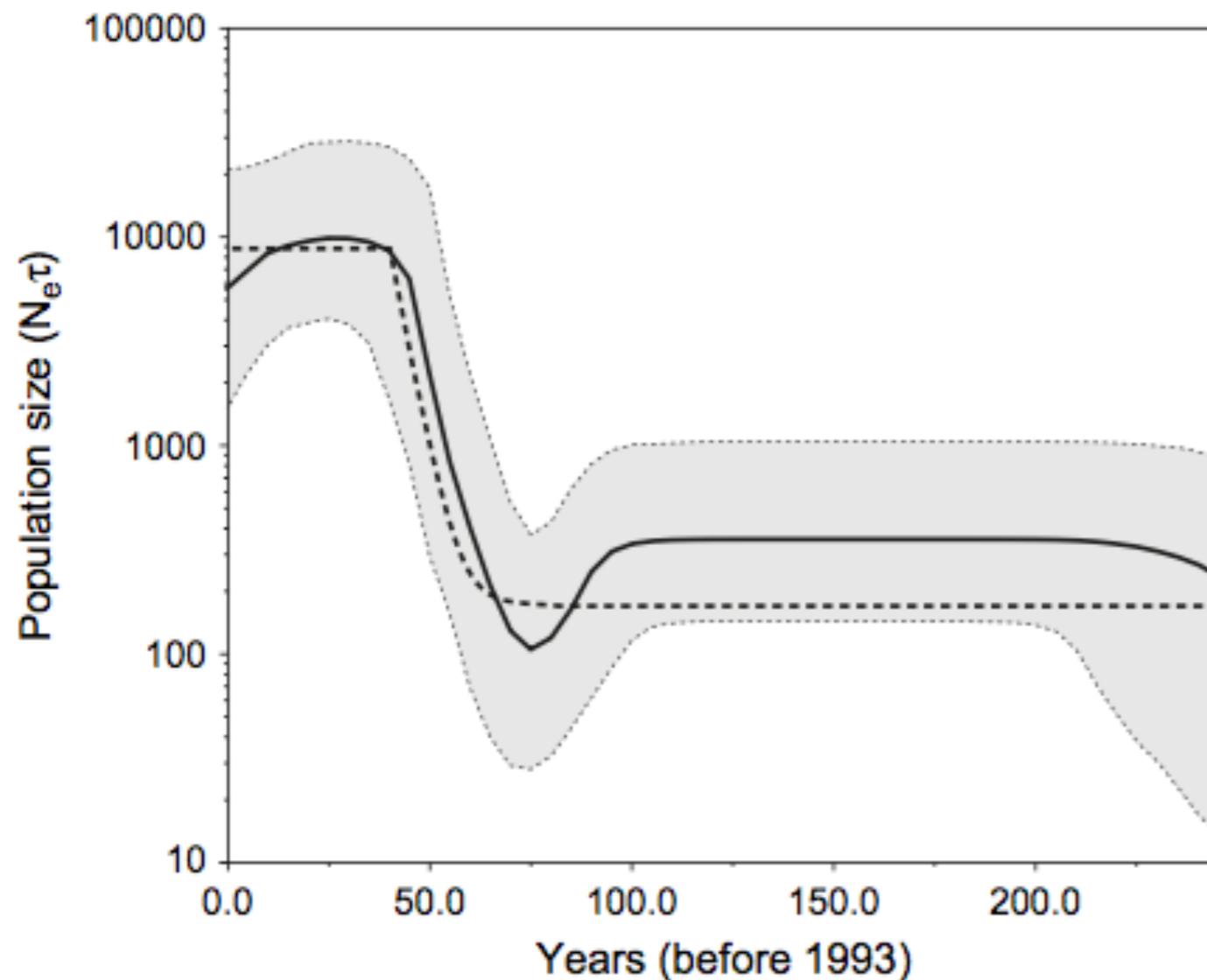
# HIV-1 “Patient 0”



- So-called “Patient 0” could not be the index case for the US HIV-1 epidemic
- US epidemic emerged from a pre-existing Caribbean epidemic ~1970

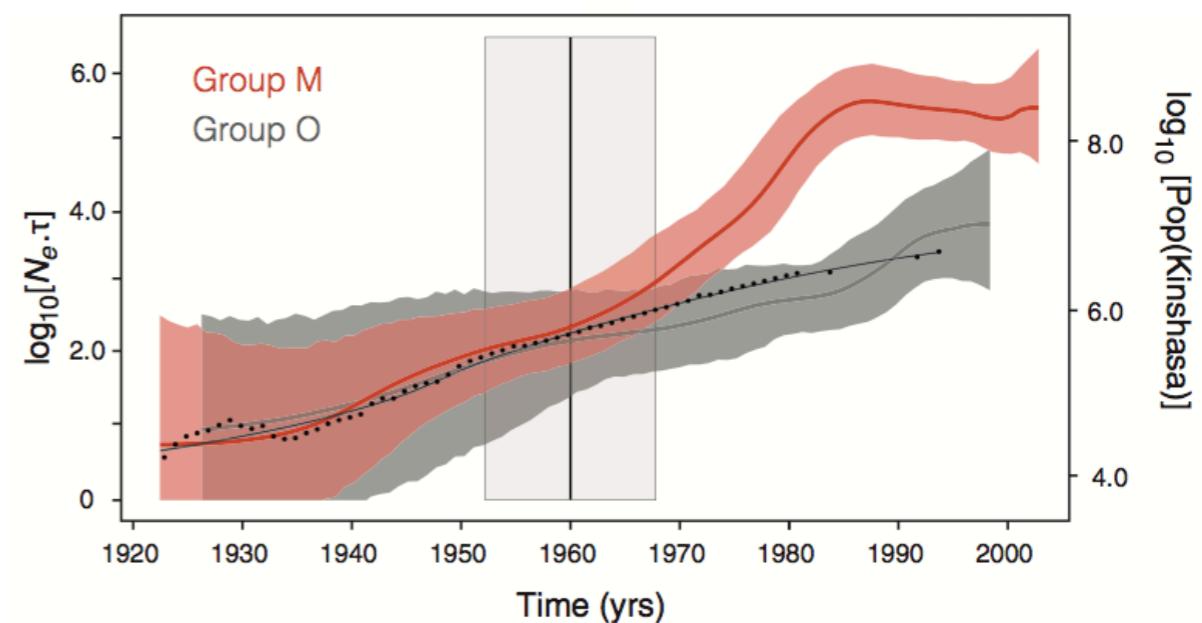
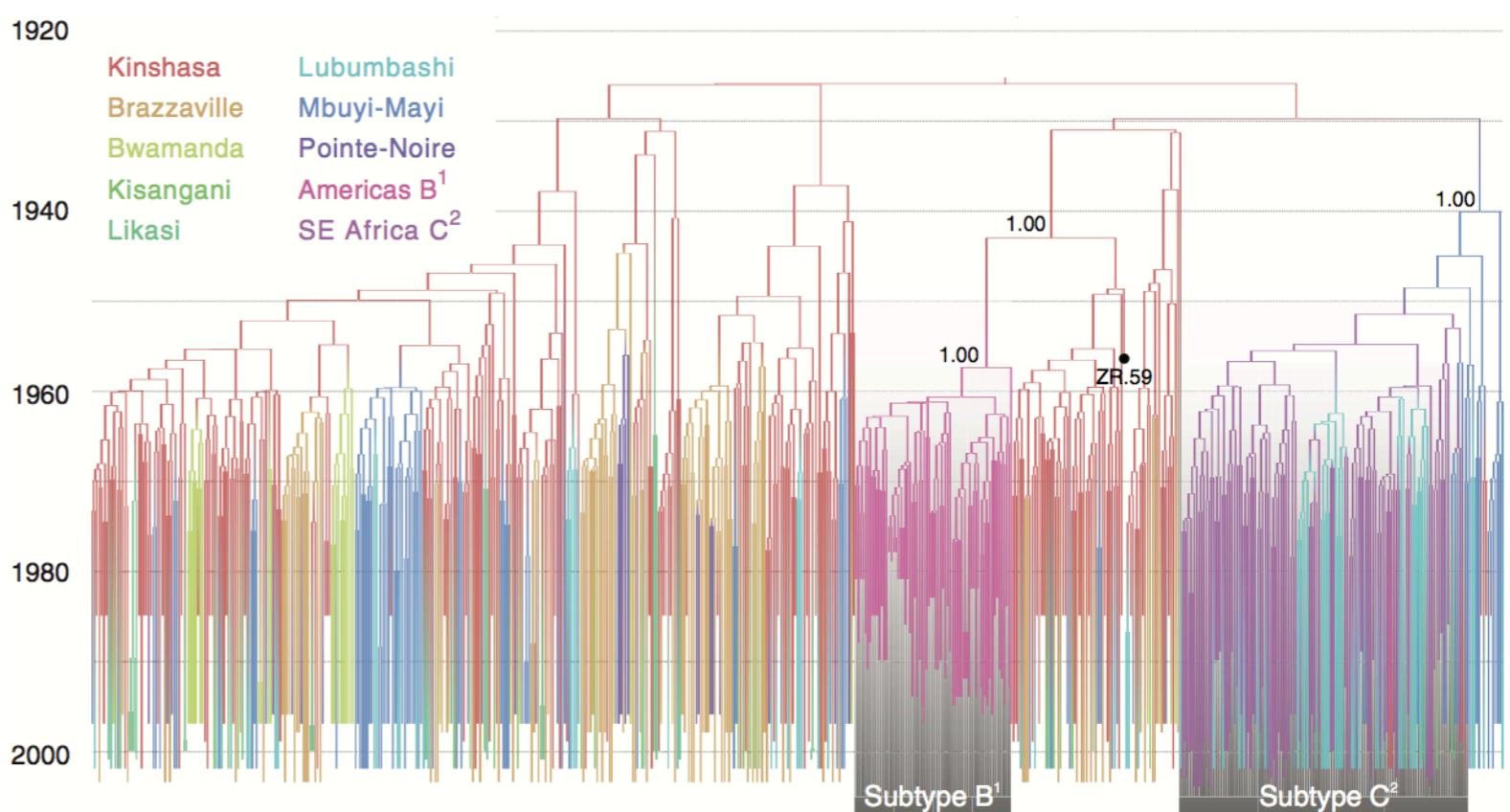
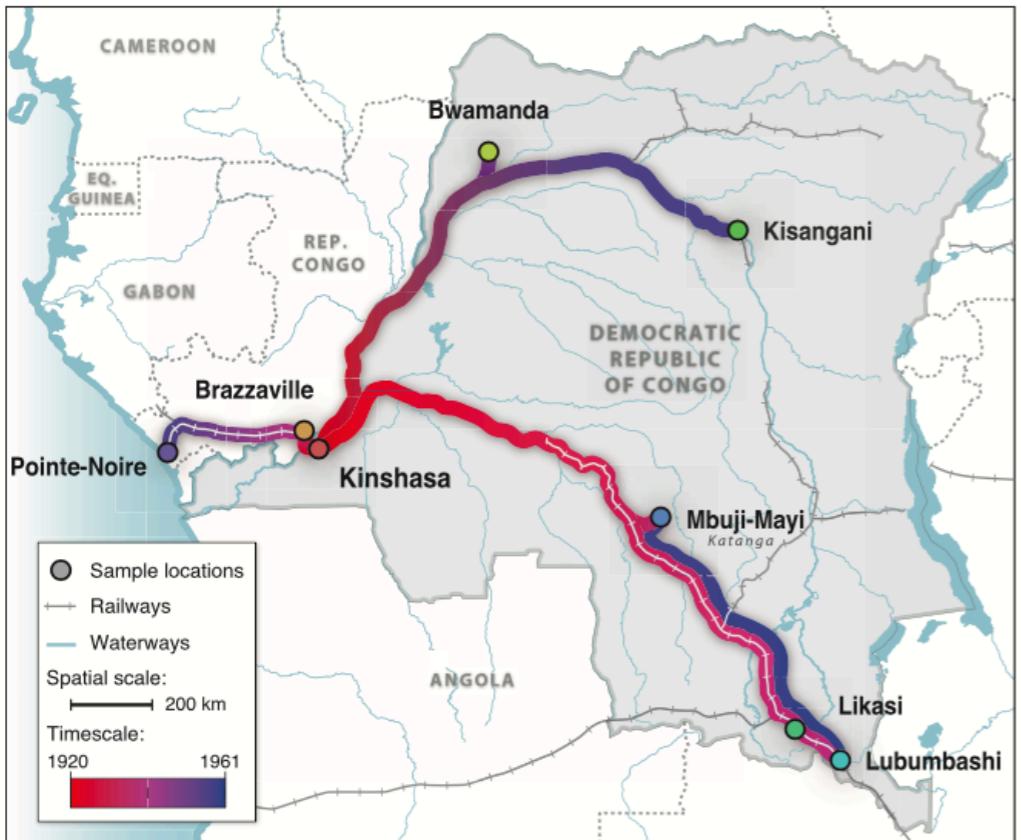
# Hepatitis C virus in Egypt

Combine with model of population dynamics

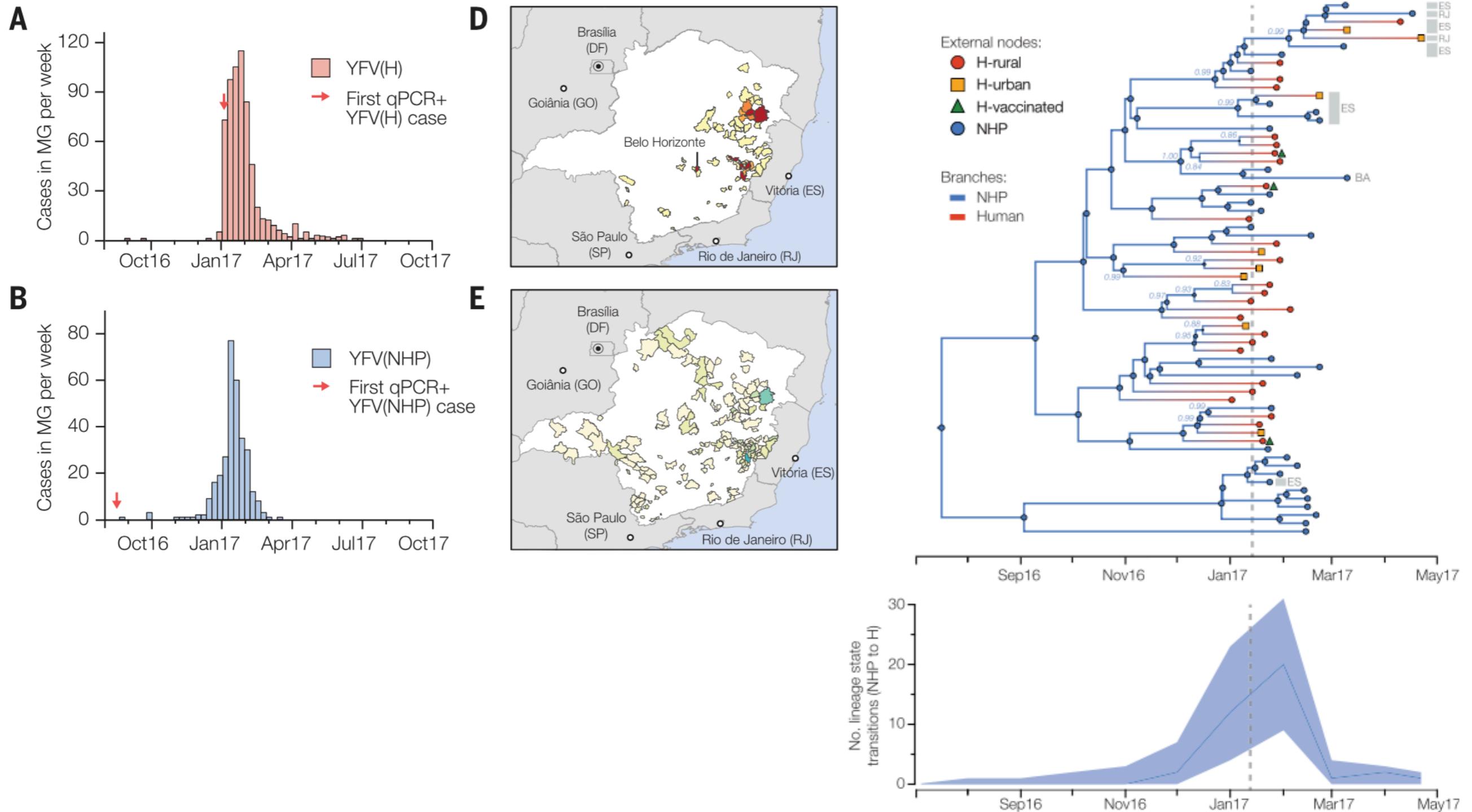


- Egypt has the highest prevalence (much higher than neighbouring countries)
- Exponential growth period during the middle of the 20<sup>th</sup> century
- Coincides with a period of possibly contaminated anti-schistosomiasis injections

# HIV-1 pandemic

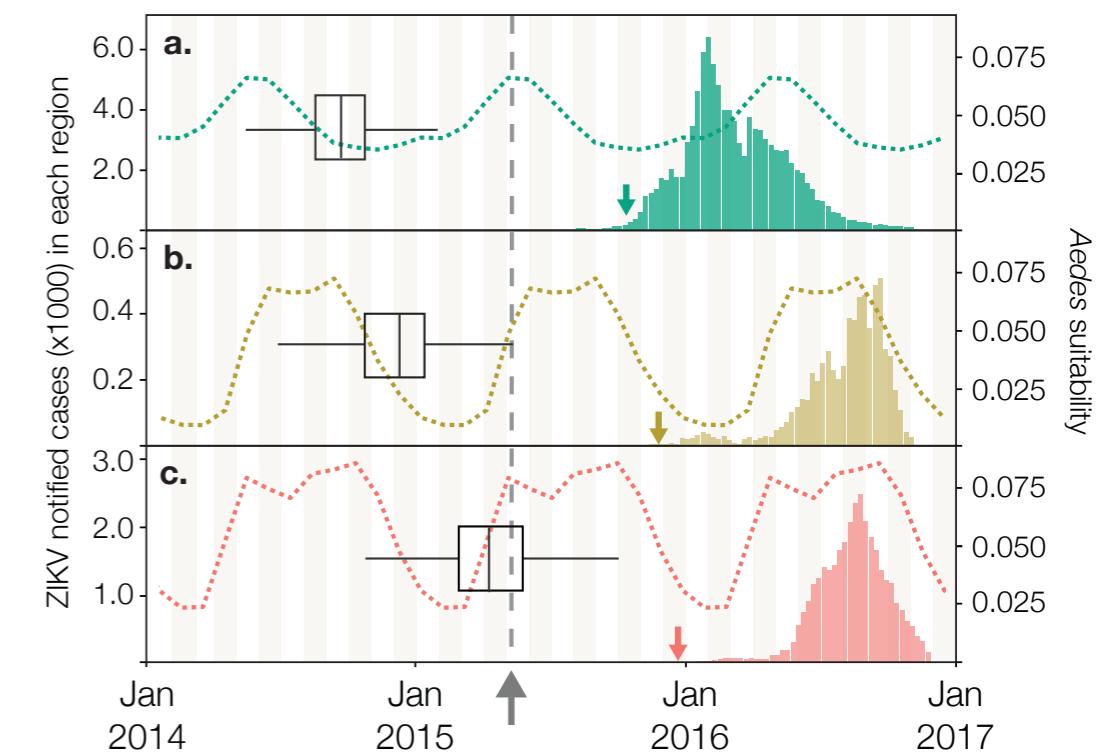
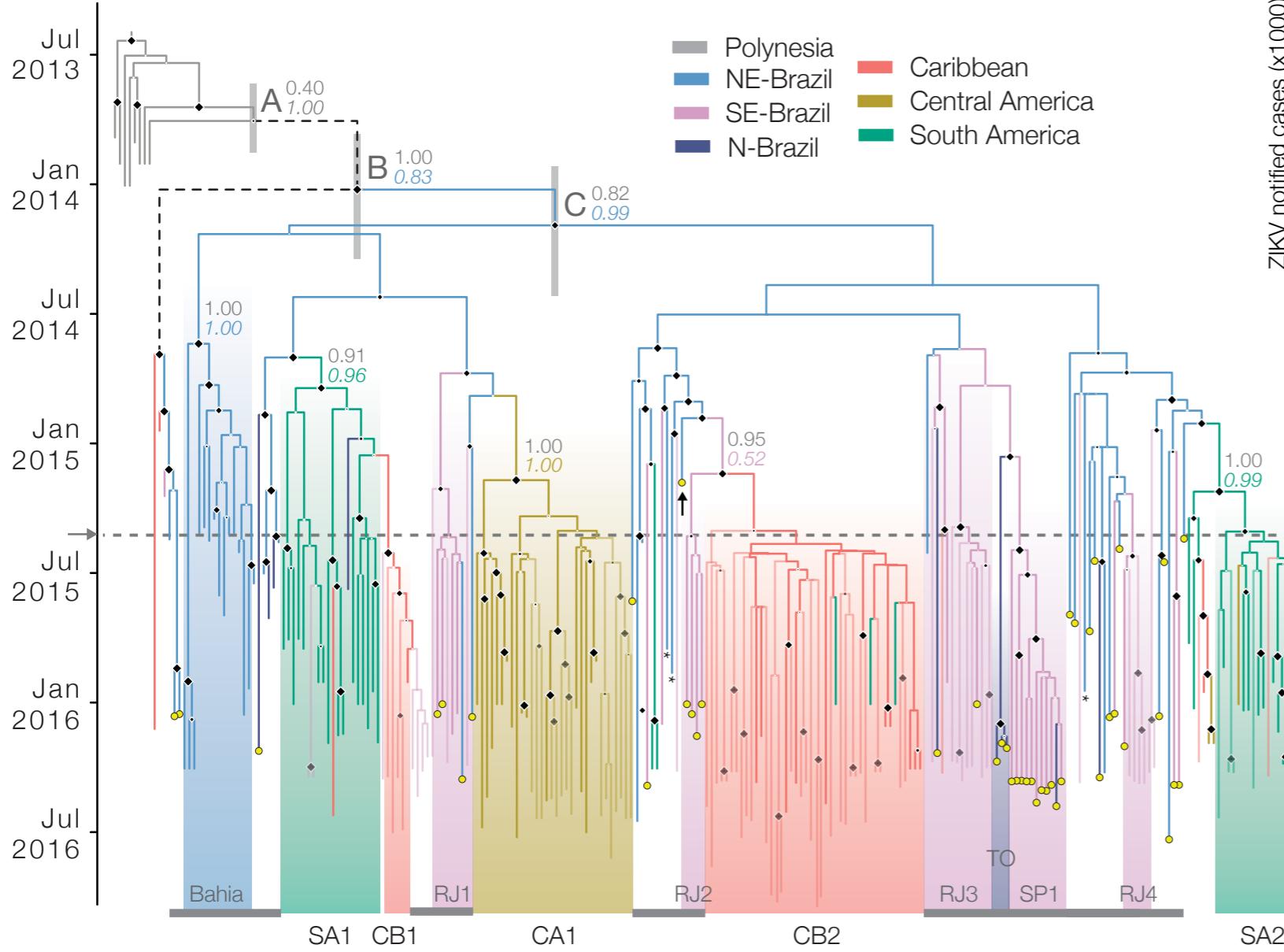


# Yellow fever in Brazil



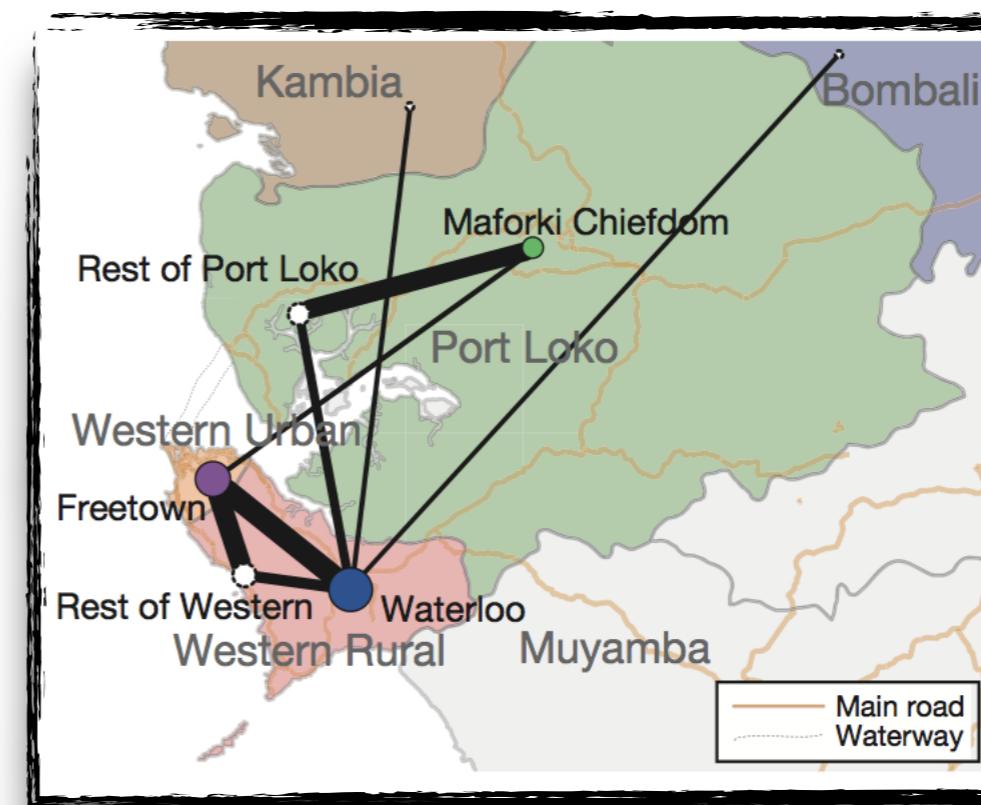
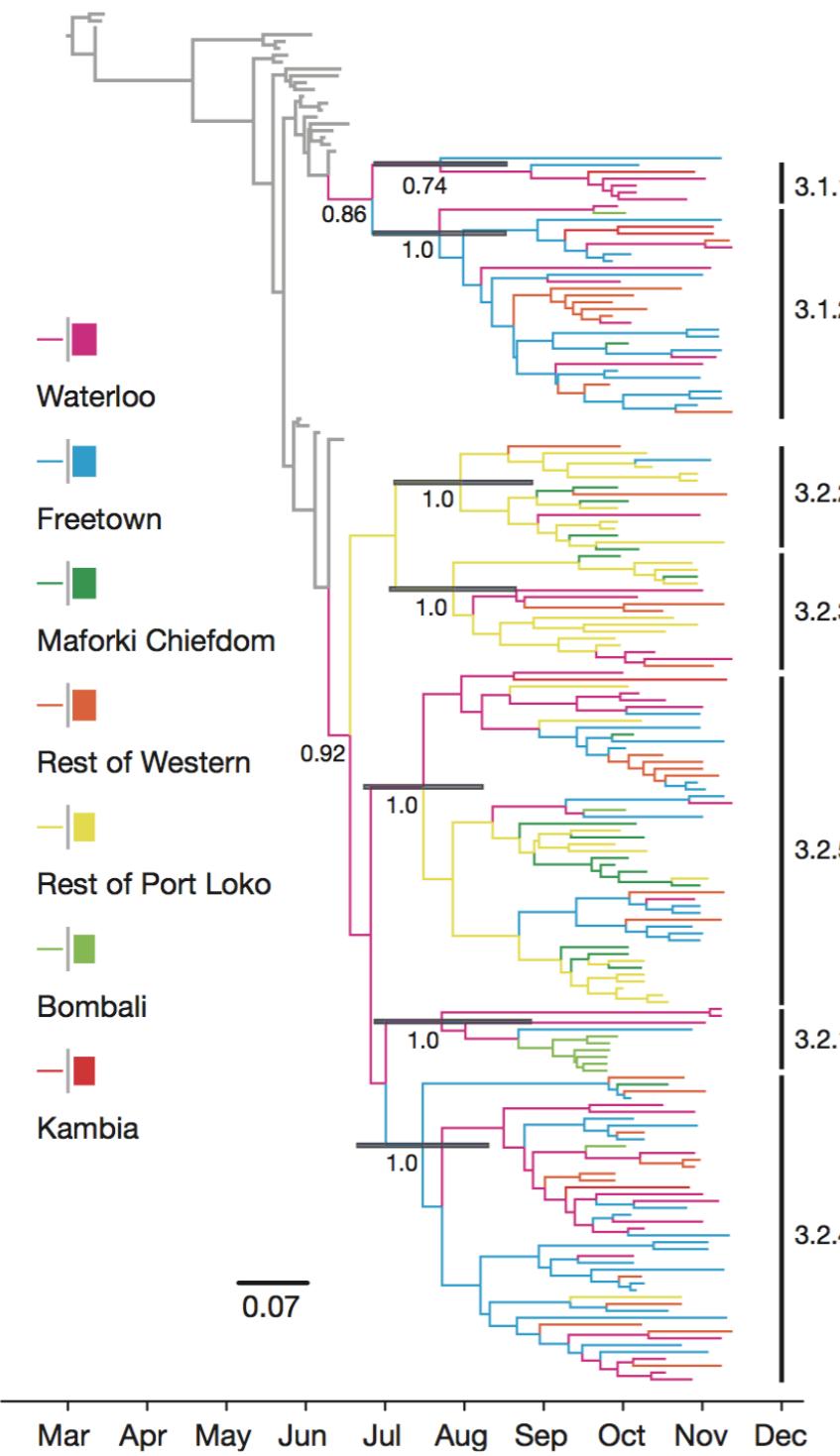
# Zika virus in the Americas

Combine with phylogeographic migration model



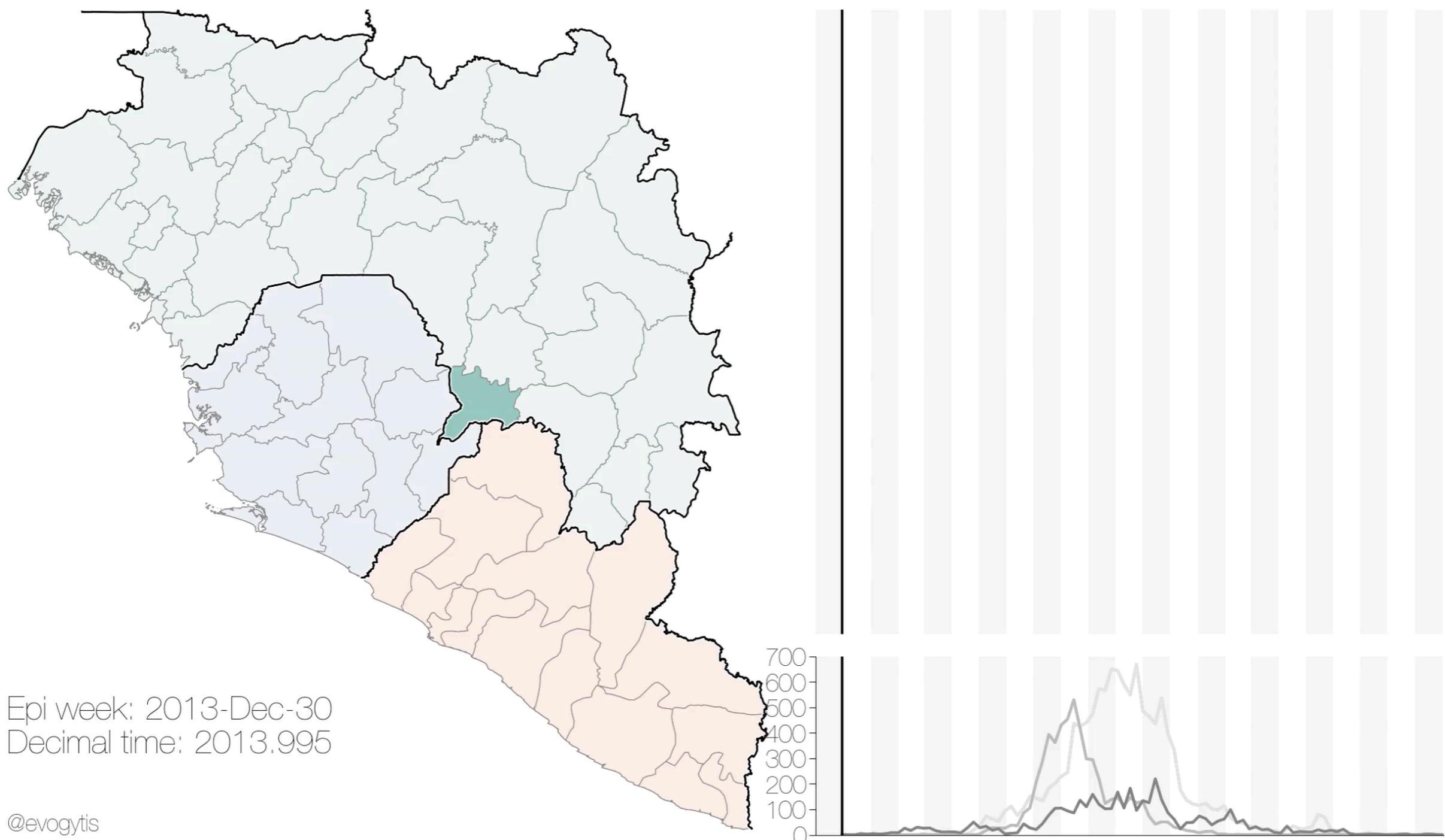
- Single introduction from French-Polynesia to Brazil
- ~1 year of cryptic transmission
- NE-Brazil is the source for the whole American epidemic

# Phylogeography: Ebola in Sierra Leone

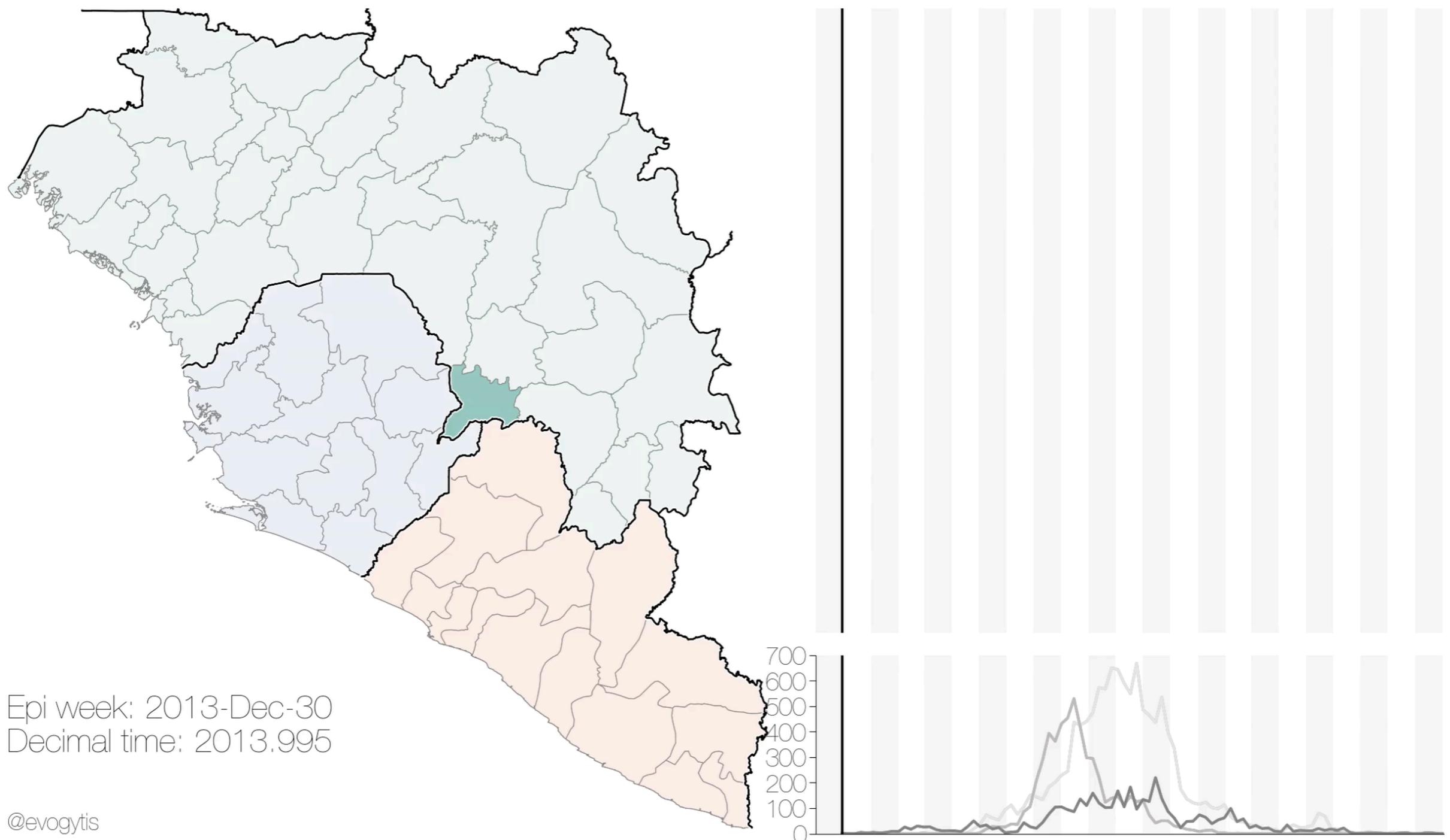


- Infer how the virus spread across the region
- Infer the most probable ancestral states
- Infer migration rates between different regions

# Phylogeography: Ebola in West Africa



# Phylogeography: Ebola in West Africa



# Seasonal Influenza (2010)

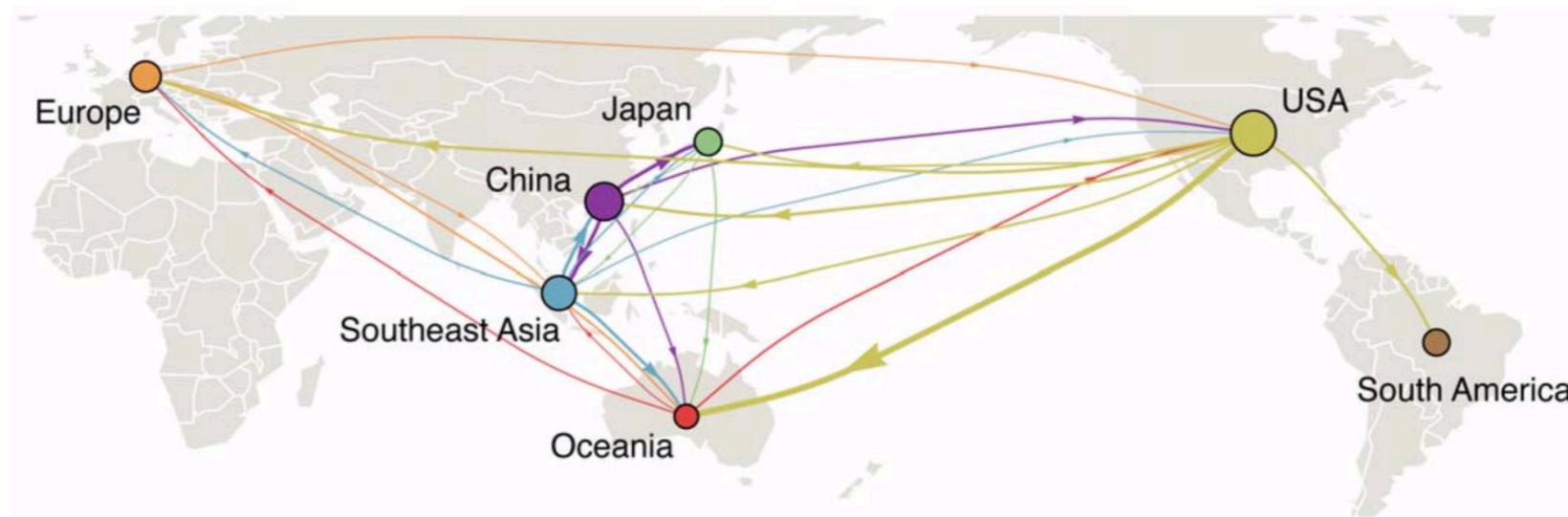
OPEN  ACCESS Freely available online

PLOS PATHOGENS

## Global Migration Dynamics Underlie Evolution and Persistence of Human Influenza A (H3N2)

Trevor Bedford<sup>1,2\*</sup>, Sarah Cobey<sup>1,2</sup>, Peter Beerli<sup>3</sup>, Mercedes Pascual<sup>1,2</sup>

**1** Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor, Michigan, United States of America, **2** Howard Hughes Medical Institute, University of Michigan, Ann Arbor, Michigan, United States of America, **3** Department of Scientific Computing, Florida State University, Tallahassee, Florida, United States of America



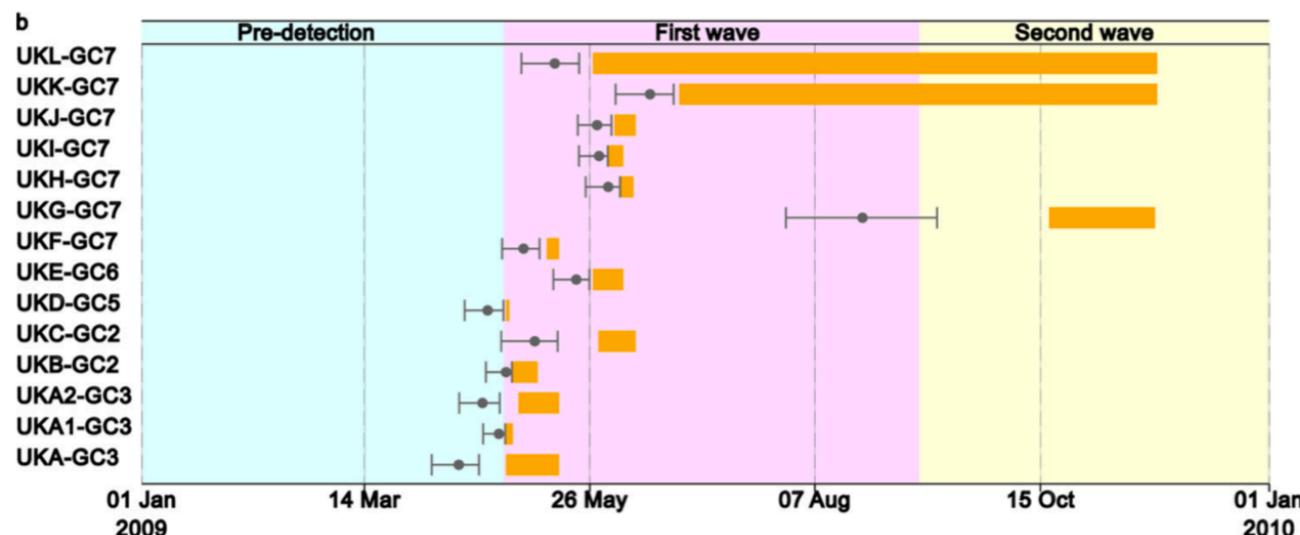
"Examining the influenza genealogy, it is apparent that regional outbreaks often result from **very few immigration events**, consistent with previous results. For example, the 2003 epidemic in Oceania appears almost completely monophyletic and can **trace its history to a single migration event** (or perhaps multiple migration events of identical strains) in early 2003."

# Pandemic H1N1 (2012)



## Evolutionary Dynamics of Local Pandemic H1N1/2009 Influenza Virus Lineages Revealed by Whole-Genome Analysis

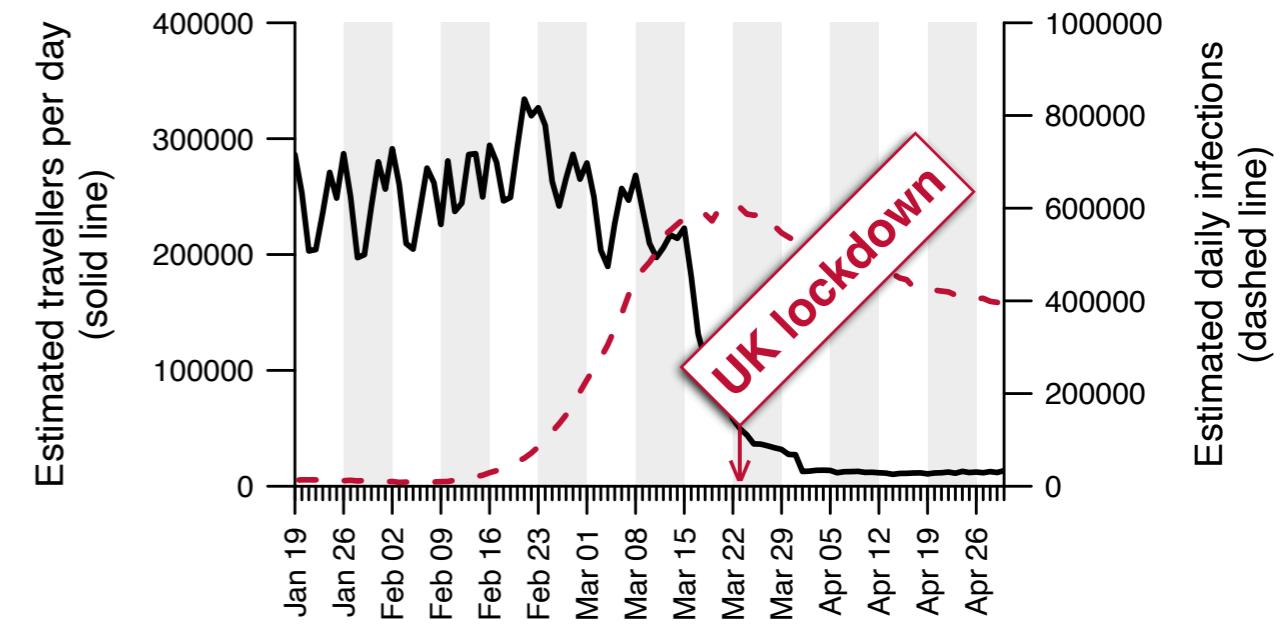
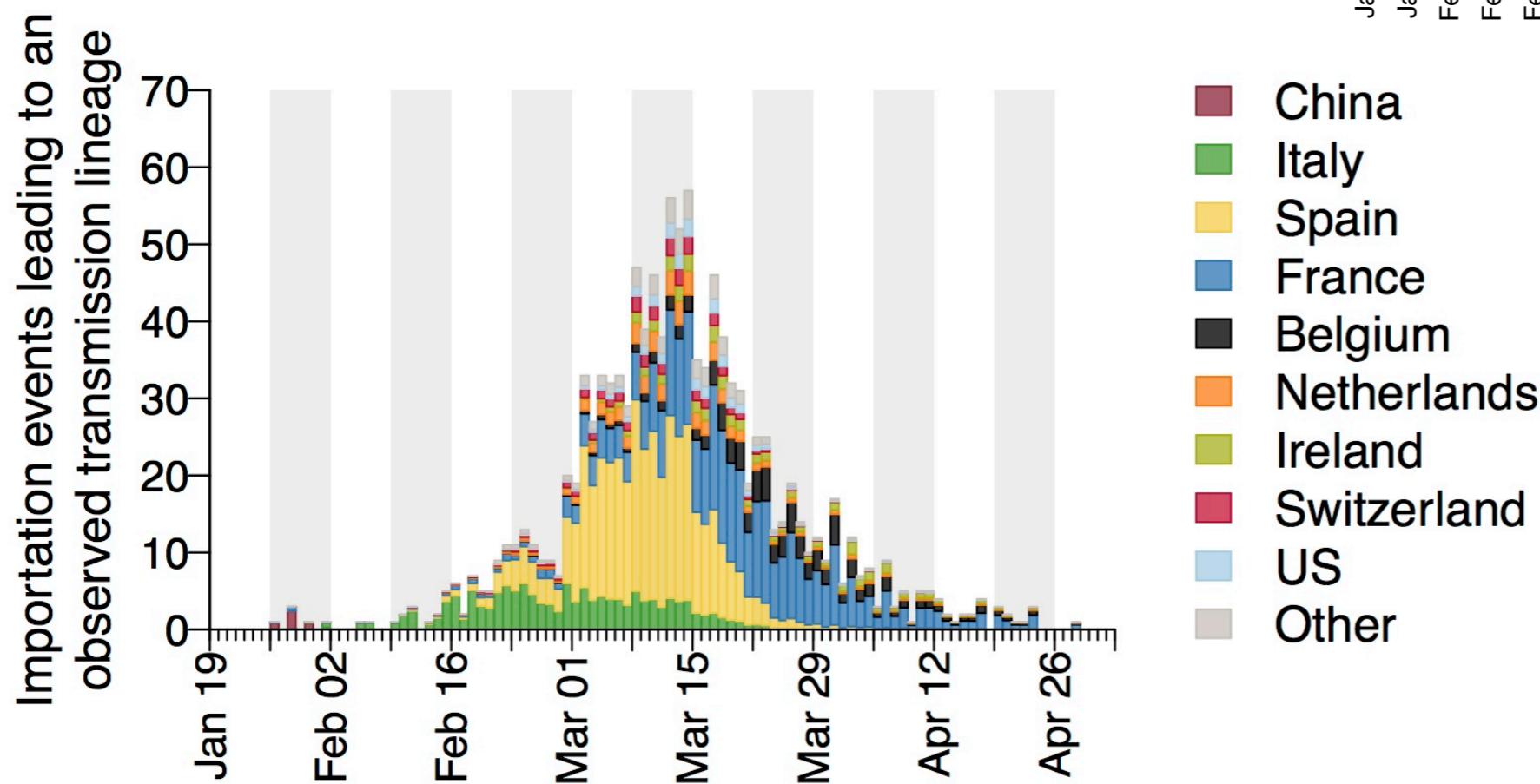
Gregory J. Baillie,<sup>a</sup> Monica Galiano,<sup>b</sup> Paul-Michael Agapow,<sup>b</sup> Richard Myers,<sup>b</sup> Rachael Chiam,<sup>a</sup> Astrid Gall,<sup>a</sup> Anne L. Palser,<sup>a</sup> Simon J. Watson,<sup>a</sup> Jessica Hedge,<sup>c</sup> Anthony Underwood,<sup>b</sup> Steven Platt,<sup>b</sup> Estelle McLean,<sup>b</sup> Richard G. Pebody,<sup>b</sup> Andrew Rambaut,<sup>c,d</sup> Jonathan Green,<sup>b</sup> Rod Daniels,<sup>e</sup> Oliver G. Pybus,<sup>f</sup> Paul Kellam,<sup>a</sup> and Maria Zambon<sup>b</sup>



"This suggests that our sampling of the lineages that arrived in the United Kingdom at the start of the outbreak was **reasonably complete** and that the **13 clusters and 52 dispersed lineages** account for much of the United Kingdom pandemic H1N1/09 virus genetic diversity."

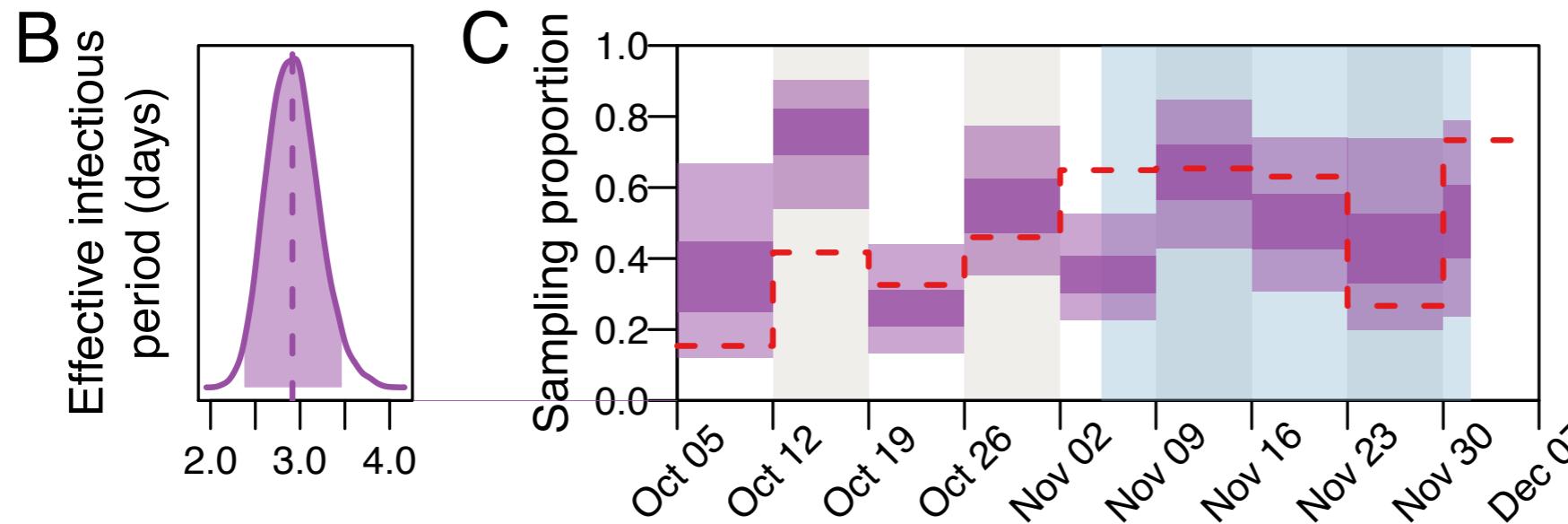
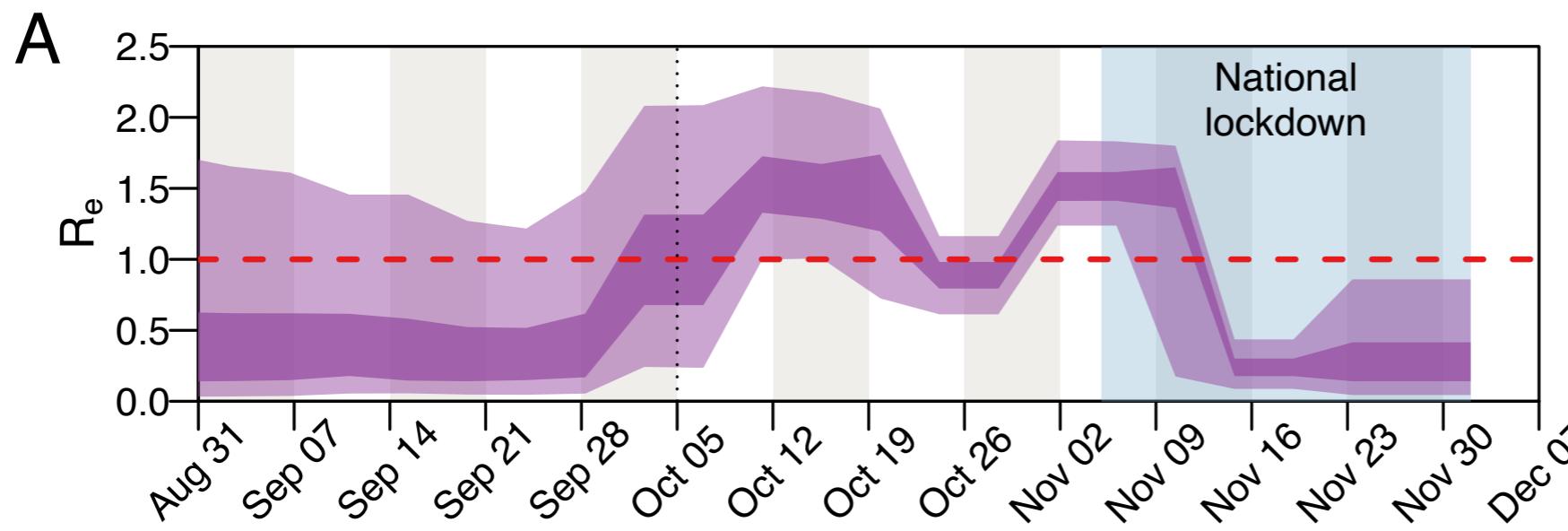
# Introduction of SARS-CoV-2 into the UK

- Uneven case ascertainment genome sampling between countries at the start of the pandemic
- Combine incidences and travel volumes to estimate importations into the UK



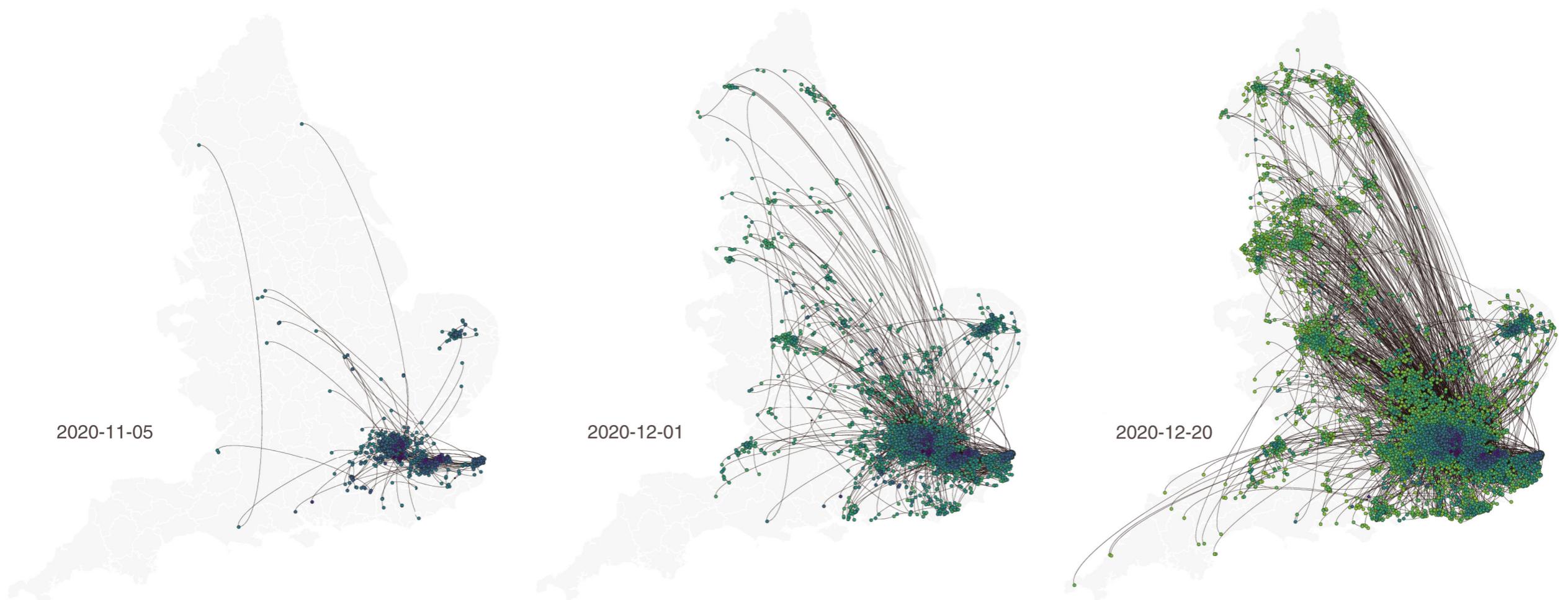
# $R_e$ of SARS-CoV-2 on a university campus (Michaelmas term 2020)

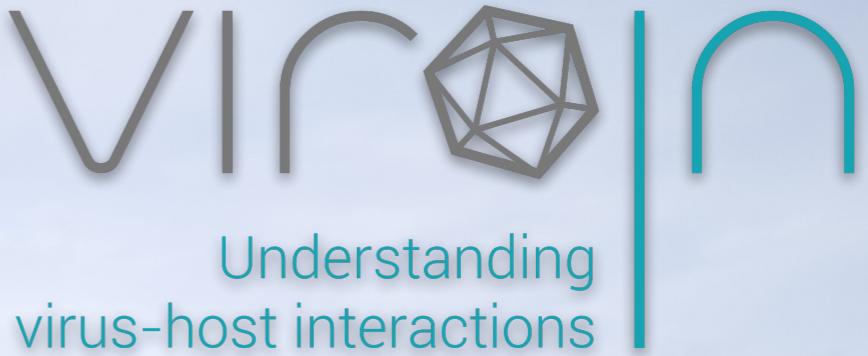
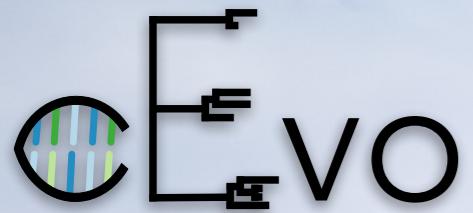
- Opt-in asymptomatic testing program for all Cambridge students
- Were campus-wide measures effective at stopping spread
- Did the national lockdown have an effect?



# Spread of Alpha VOC in England

- A new variant (B.1.1.7) appeared in southeast England in September 2020 and spread rapidly
- Later designated as the Alpha variant of concern (VOC)
- Use a phylogeographic model to track spread





# Bayesian inference recap

Louis du Plessis

# On the program for today

---

(1) What is phylodynamics?

**(2) Bayesian inference recap**

(3) BEAST2 introduction

**Tutorial:** Molecular clock dating (part i)

(4) Molecular clock models

**Tutorial:** Molecular clock dating (part ii)

(5) Setting priors

**Tutorial:** Phylodynamics (part i)

(6) Tree priors

**Tutorial:** Phylodynamics (part ii)

## **data**

- Samples drawn from a realisation of some **stochastic process**
- Typically an alignment of DNA or RNA sequences
- Can also be amino acids or codons

## **model**

- **Model** is a mathematical description of the **process** that generated the data
- Each model has a number of **model parameters**
- Parameters are **random variables**

## **hypothesis**

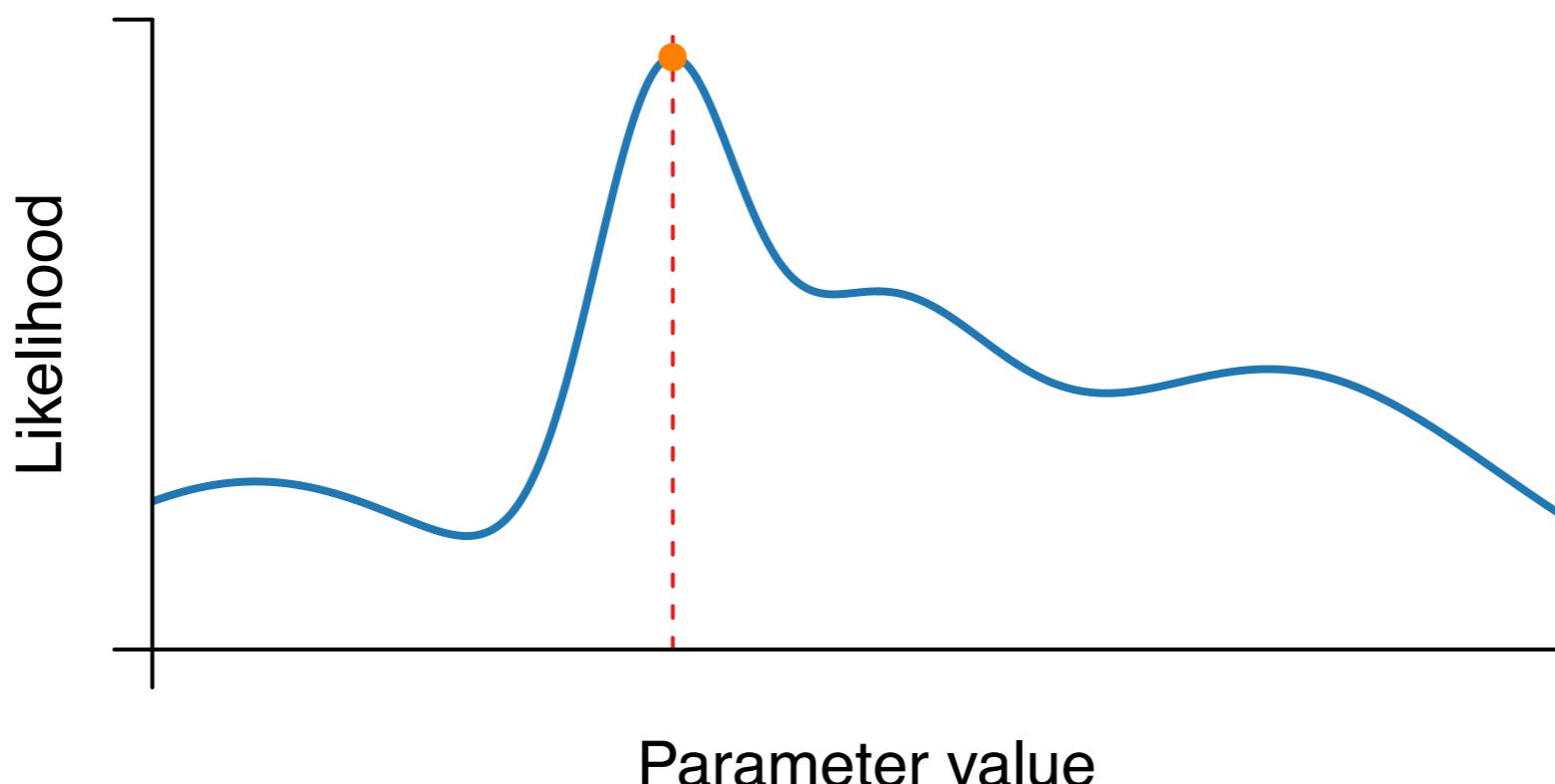
- **Hypothesis** is some statement about the model parameters
- Hypothesis assigns **values** to the model parameters
- Data used to decide if hypothesis is reasonable or not
- May only concern some model parameters. The rest are **nuisance parameters**

# Maximum-likelihood inference

---

**Likelihood** →  $P(\text{data} \mid \text{hypothesis, model})$

- Likelihood is proportional to the probability of observing the data given a **hypothesis**
- Find the model parameters that maximise the likelihood



# Inference frameworks

---

## Frequentist

- Probabilities refer to the outcome of experiments
- Probability is the frequency of outcomes if an experiment is repeated many times
- Likelihood is the degree to which data supports a hypothesis

## Bayesian

- Probability is a measure of the plausibility of propositions conditional on available information
- Data and model parameters are described by probabilities
- Probability reflects our belief in a hypothesis as well as representing the outcome of experiments
- Hypotheses have probabilities in the absence of data

# DID THE SUN JUST EXPLODE? (IT'S NIGHT, SO WE'RE NOT SURE.)

THIS NEUTRINO DETECTOR MEASURES WHETHER THE SUN HAS GONE NOVA.

THEN, IT ROLLS TWO DICE. IF THEY BOTH COME UP SIX, IT LIES TO US. OTHERWISE, IT TELLS THE TRUTH.

LET'S TRY.

DETECTOR! HAS THE SUN GONE NOVA?

ROLL

YES.



FREQUENTIST STATISTICIAN:

THE PROBABILITY OF THIS RESULT HAPPENING BY CHANCE IS  $\frac{1}{36} = 0.027$ .

SINCE  $p < 0.05$ , I CONCLUDE THAT THE SUN HAS EXPLODED.



BAYESIAN STATISTICIAN:

BET YOU \$50 IT HASN'T.



# Bayesian inference

---

## Prior → $P(\text{hypothesis} \mid \text{model})$

- Have some degree of belief in our hypothesis
- From external sources, previous analyses etc.

## Likelihood → $P(\text{data} \mid \text{hypothesis, model})$

- Likelihood is proportional to the probability of observing the data given a hypothesis

## Posterior → $P(\text{hypothesis} \mid \text{data, model})$

- Combines information from the data (**likelihood**) and previous knowledge (**prior**)

# Bayesian inference

---

## Prior → $P(\text{hypothesis})$

- Have some degree of belief in our hypothesis
- From external sources, previous analyses etc.

## Likelihood → $P(\text{data} \mid \text{hypothesis})$

- Likelihood is proportional to the probability of observing the data given a hypothesis

## Posterior → $P(\text{hypothesis} \mid \text{data})$

- Combines information from the data (**likelihood**) and previous knowledge (**prior**)

**Bayesian inference estimates the posterior probability distribution instead of a likelihood curve**

# Bayesian inference

---

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

Likelihood

Prior

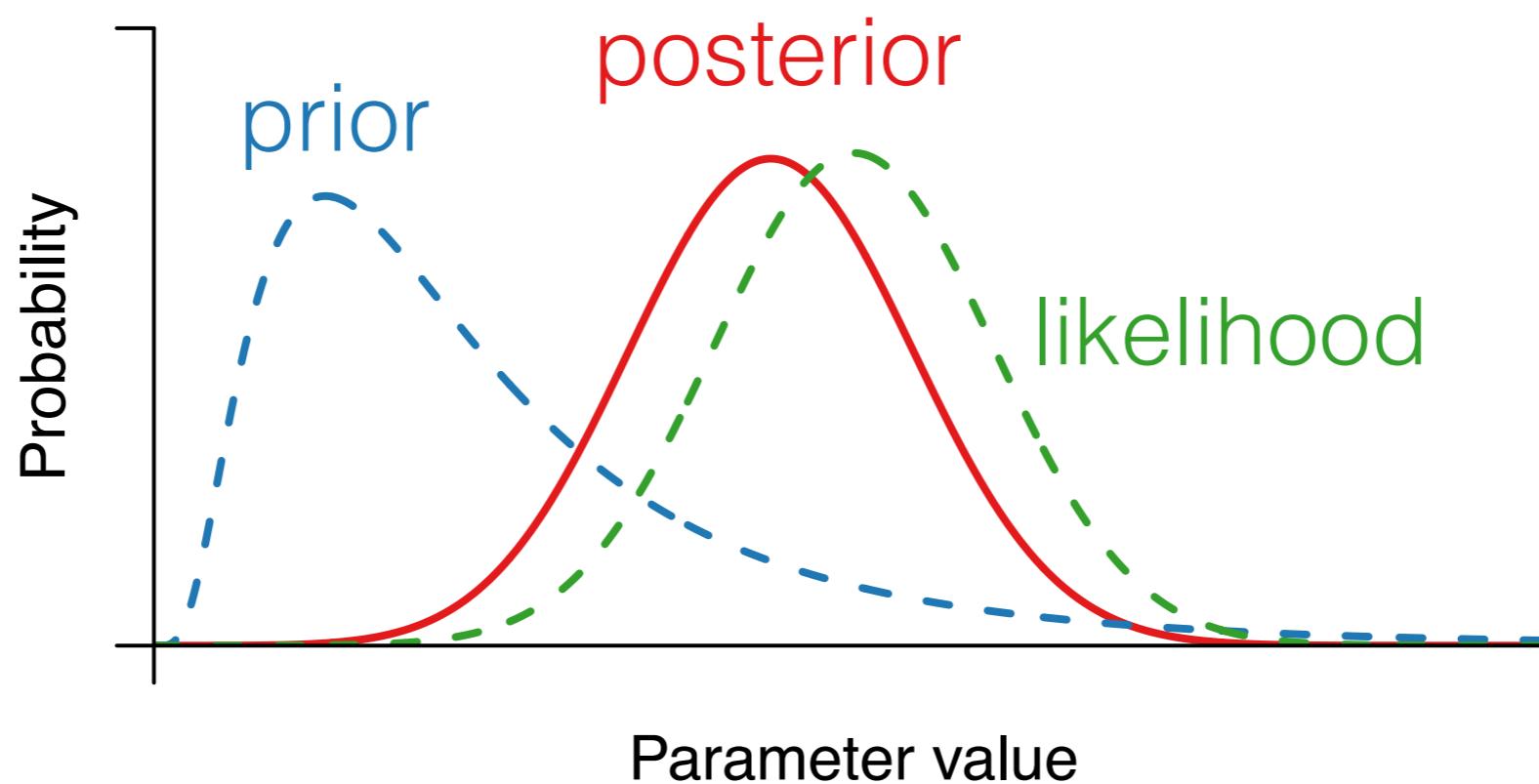
Posterior

Model  
evidence

# Bayesian inference

---

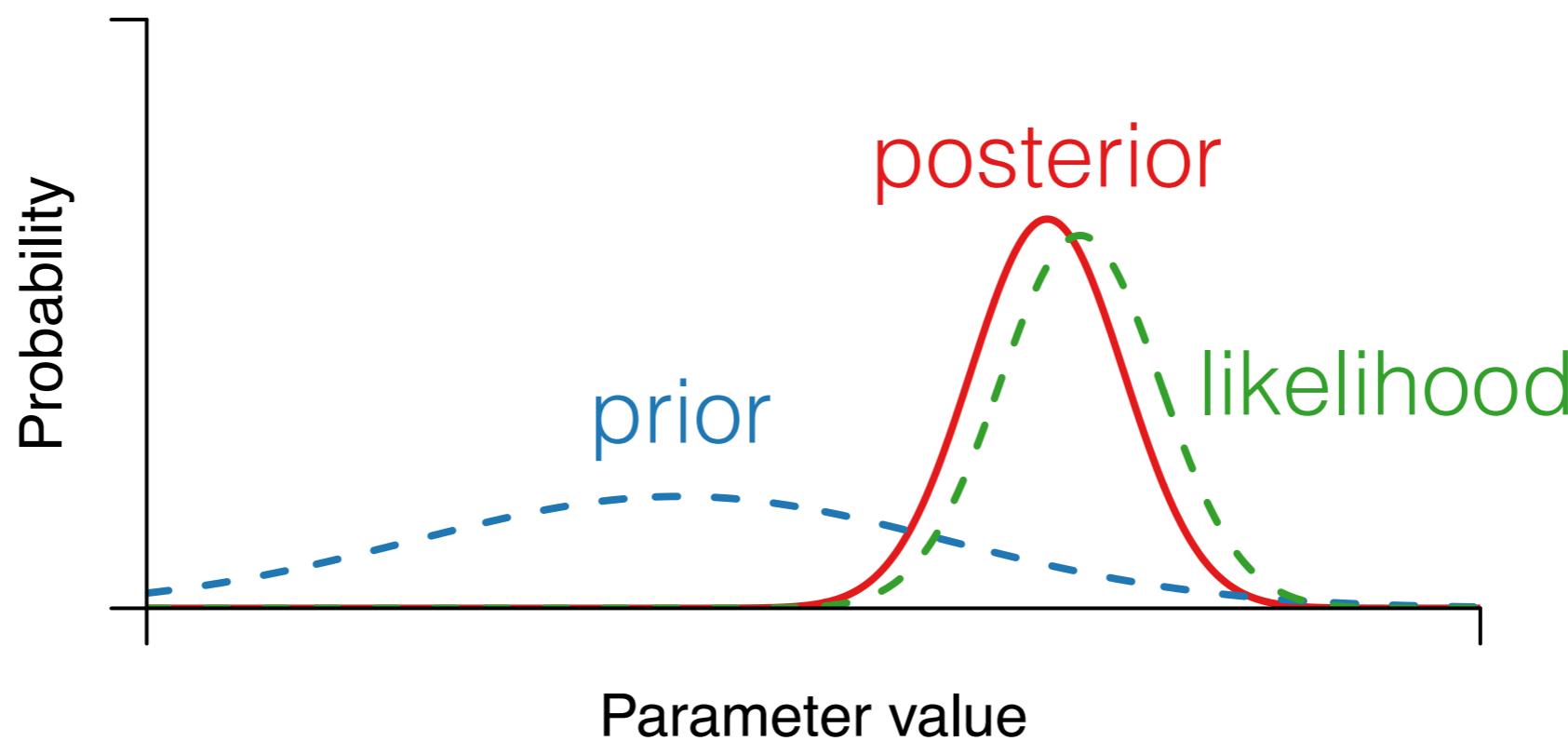
Prior distribution is combined (updated) with the likelihood to yield the **posterior distribution**



# Bayesian inference

---

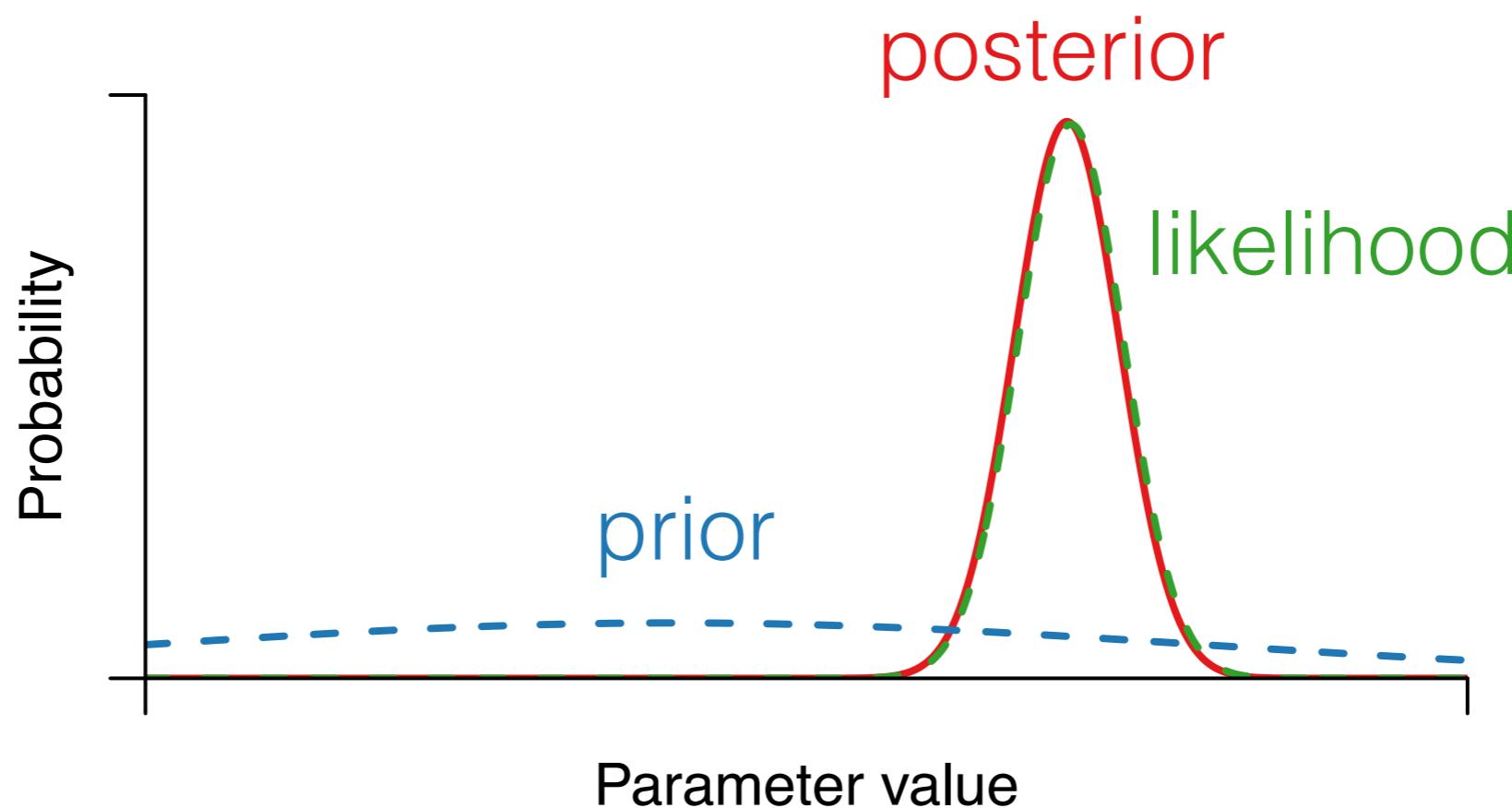
The more informative the data, the less effect the prior has



# Bayesian inference

---

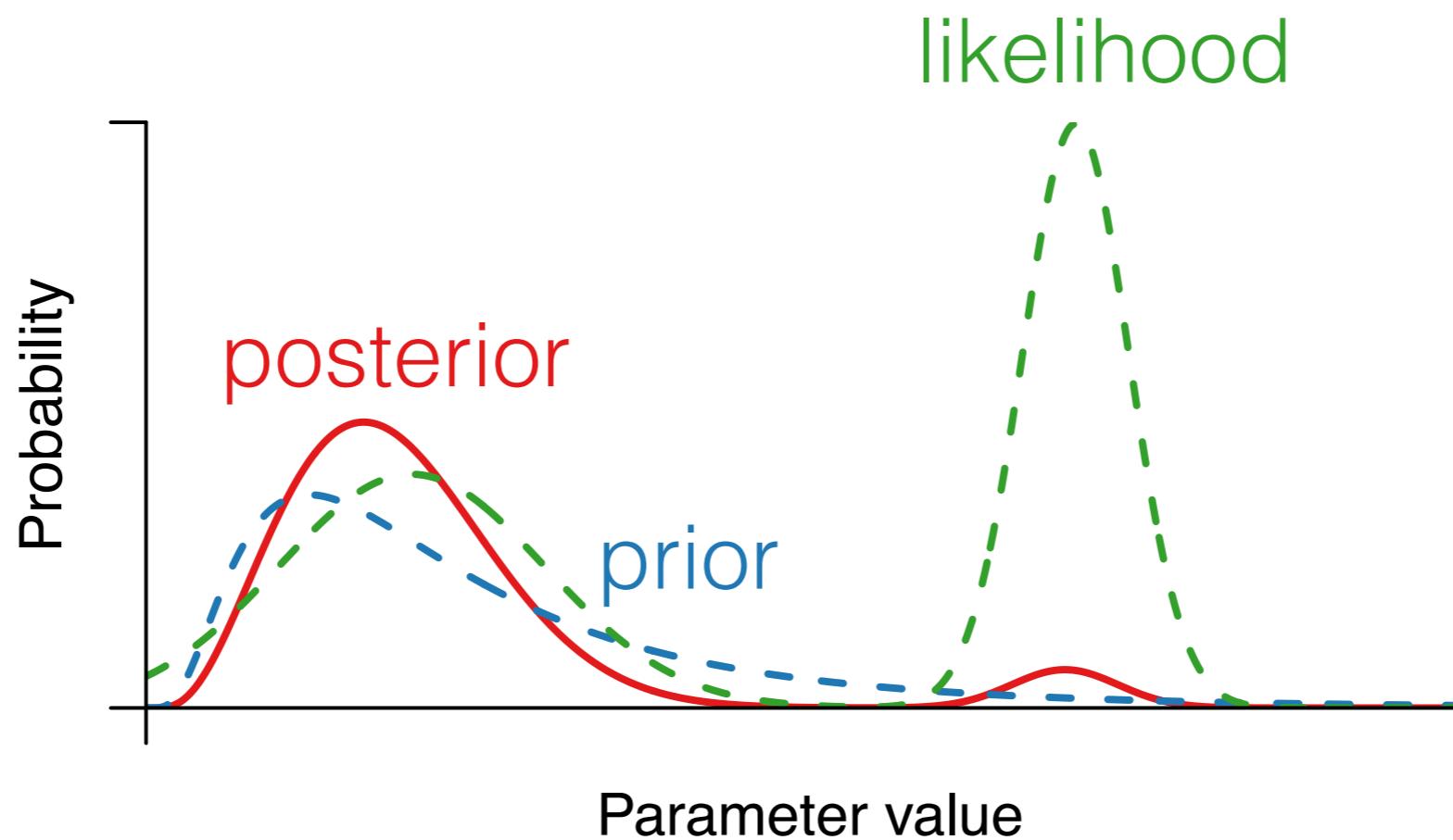
If the prior is very uninformative and the likelihood very informative the posterior will be almost exactly the likelihood



# Bayesian inference

---

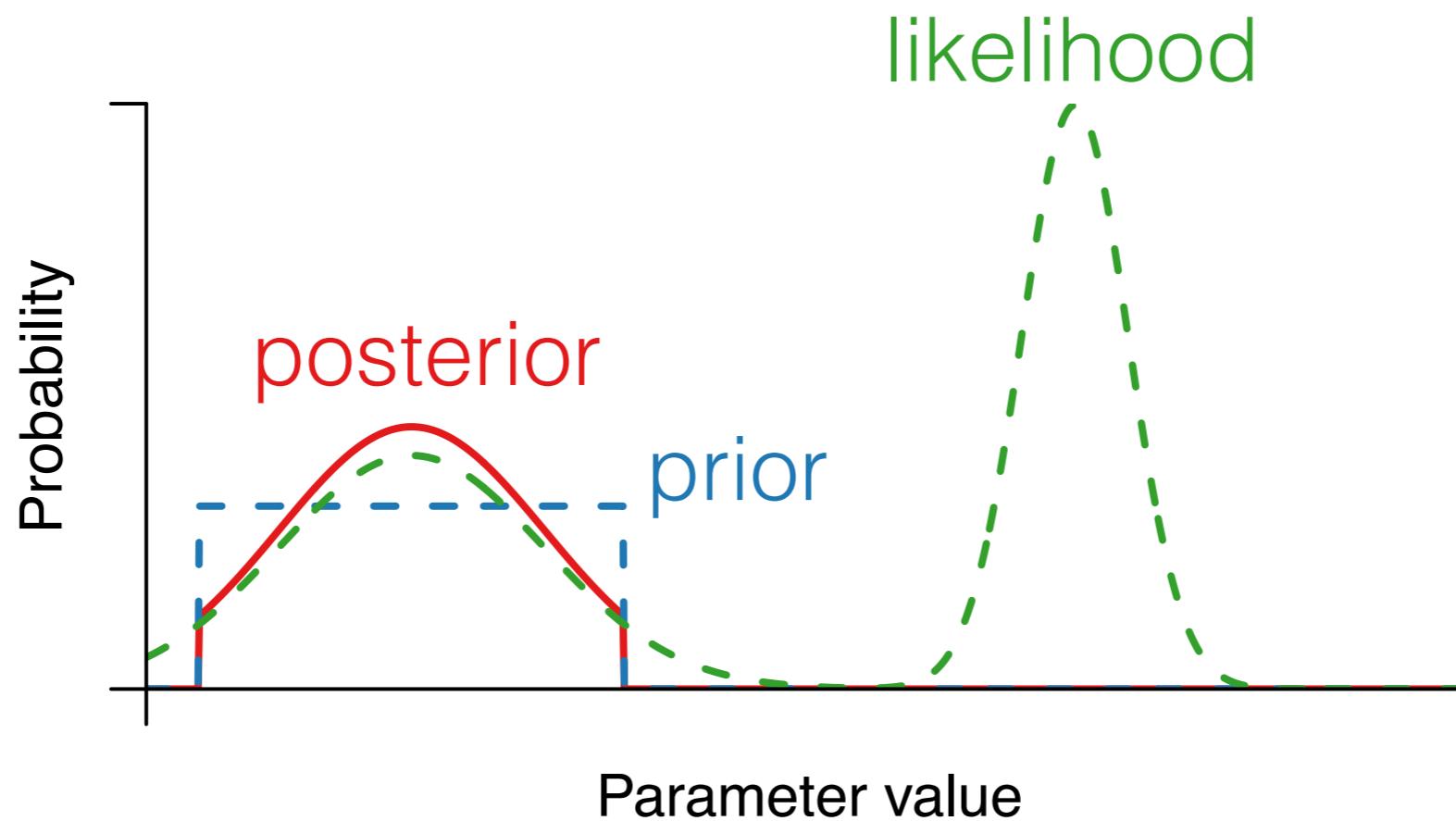
If the likelihood is multimodal the posterior is usually also multimodal even when the prior is informative



# Bayesian inference

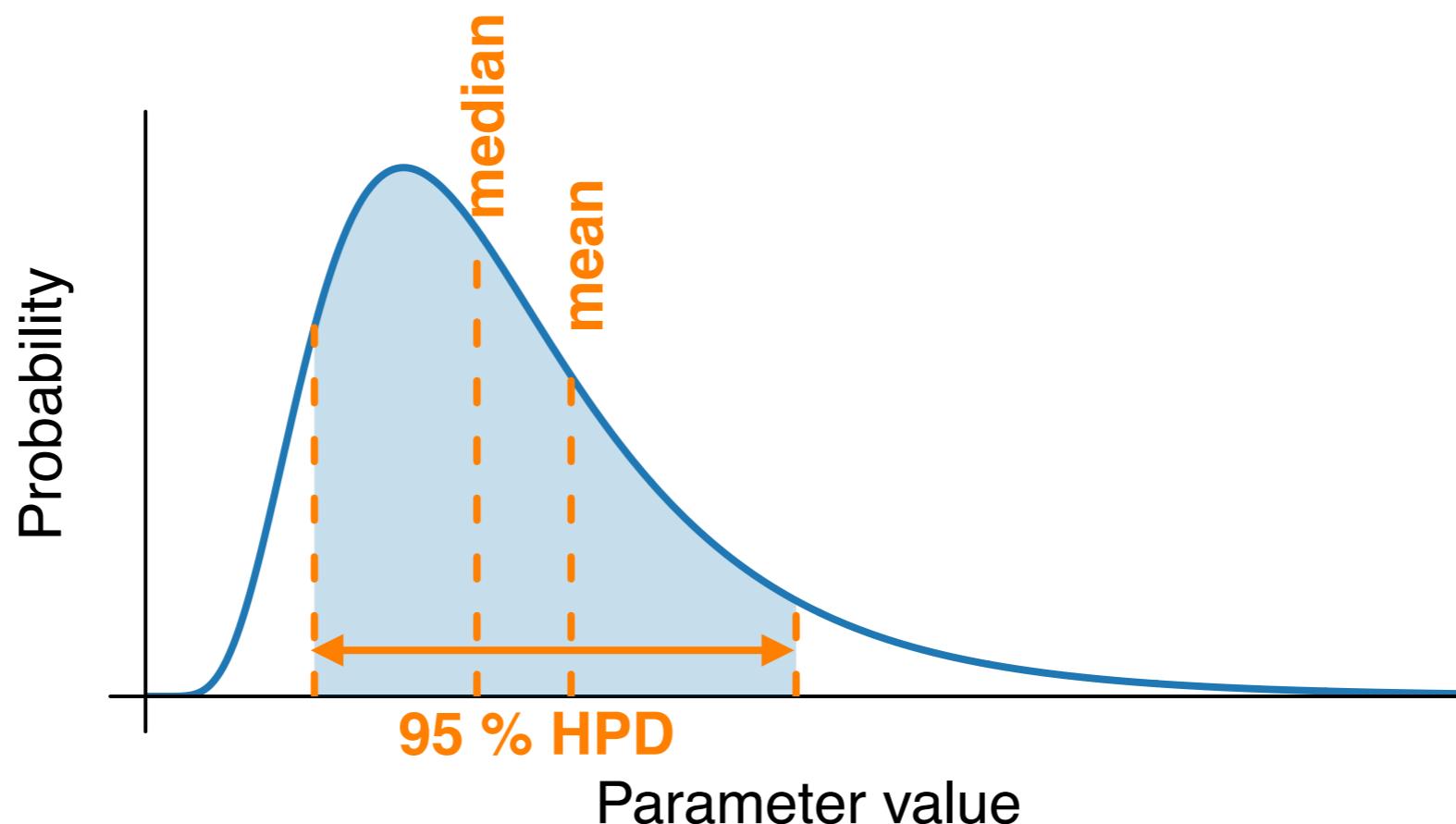
---

Can use the prior to restrict the posterior to specific ranges,  
e.g. biologically realistic values



# Summarising posteriors

---



95% HPD interval is the smallest interval that contains 95% of the posterior probability

# Bias and precision

---



**Biased and imprecise**



**Biased and precise**



**Unbiased and imprecise**



**Unbiased and precise**

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

### Prior → $P(\text{hypothesis})$

- Original probability for the model parameters/components
- Belief in our hypothesis
- All parameters have priors, whether you specify them or not!

### Likelihood → $P(\text{data} \mid \text{hypothesis})$

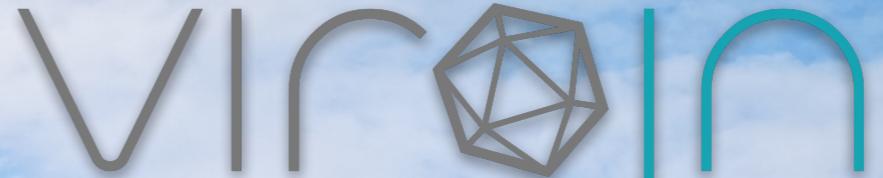
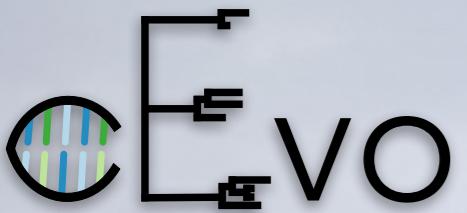
- Probability of data given parameters (defined by model)

### Posterior → $P(\text{hypothesis} \mid \text{data})$

- Updated probability for the model parameters in light of the data

### Model evidence → $P(\text{data}) = P(\text{data} \mid \text{model})$

- Probability for data given model (any combination of parameters)
- Not generally used for parameter inference
- Used for Bayesian model selection



Understanding  
virus-host interactions

101010 010010 100010 010010  
01110011



# BEAST2 introduction

Louis du Plessis

ETH zürich

DBSSE

# On the program for today

---

- (1) What is phylodynamics?
- (2) Bayesian inference recap

## **(3) BEAST2 introduction**

**Tutorial:** Molecular clock dating (part i)

- (4) Molecular clock models

**Tutorial:** Molecular clock dating (part ii)

- (5) Setting priors

**Tutorial:** Phylodynamics (part i)

- (6) Tree priors

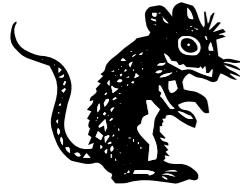
**Tutorial:** Phylodynamics (part ii)

# beast      noun

\bēst\

## Definition of *beast*

1. any nonhuman animal, especially a large, four-footed mammal
2. a contemptible person
3. something formidably difficult to control or deal with
4. **the beast**, the Antichrist. Rev. 13:18



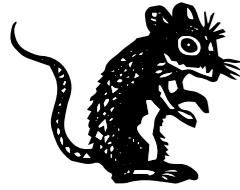
**beast2**

**noun**

\bēst-tōō\

## Definition of *beast2*

1. **B**ayesian **E**volutionary **A**nalysis by **S**ampling **T**rees **2**
2. a cross-platform package for performing Bayesian inference using MCMC

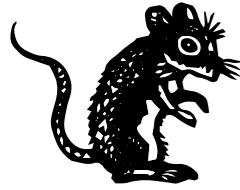


**beast2**      noun

\bēst-tōō\

## Definition of *beast2*

1. **B**ayesian **E**volutionary **A**nalysis by **S**ampling **T**rees **2**
2. a cross-platform package for performing Bayesian inference using MCMC *with emphasis on phylogenetic analysis of molecular sequences*

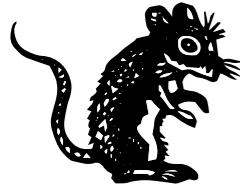


**beast2**      noun

\bēst-tōō\

## Definition of *beast2*

1. **B**ayesian **E**volutionary **A**nalysis by **S**ampling **T**rees **2**
2. a *modular, extensible, cross-platform package for performing Bayesian inference using MCMC with emphasis on phylogenetic analysis of molecular sequences*



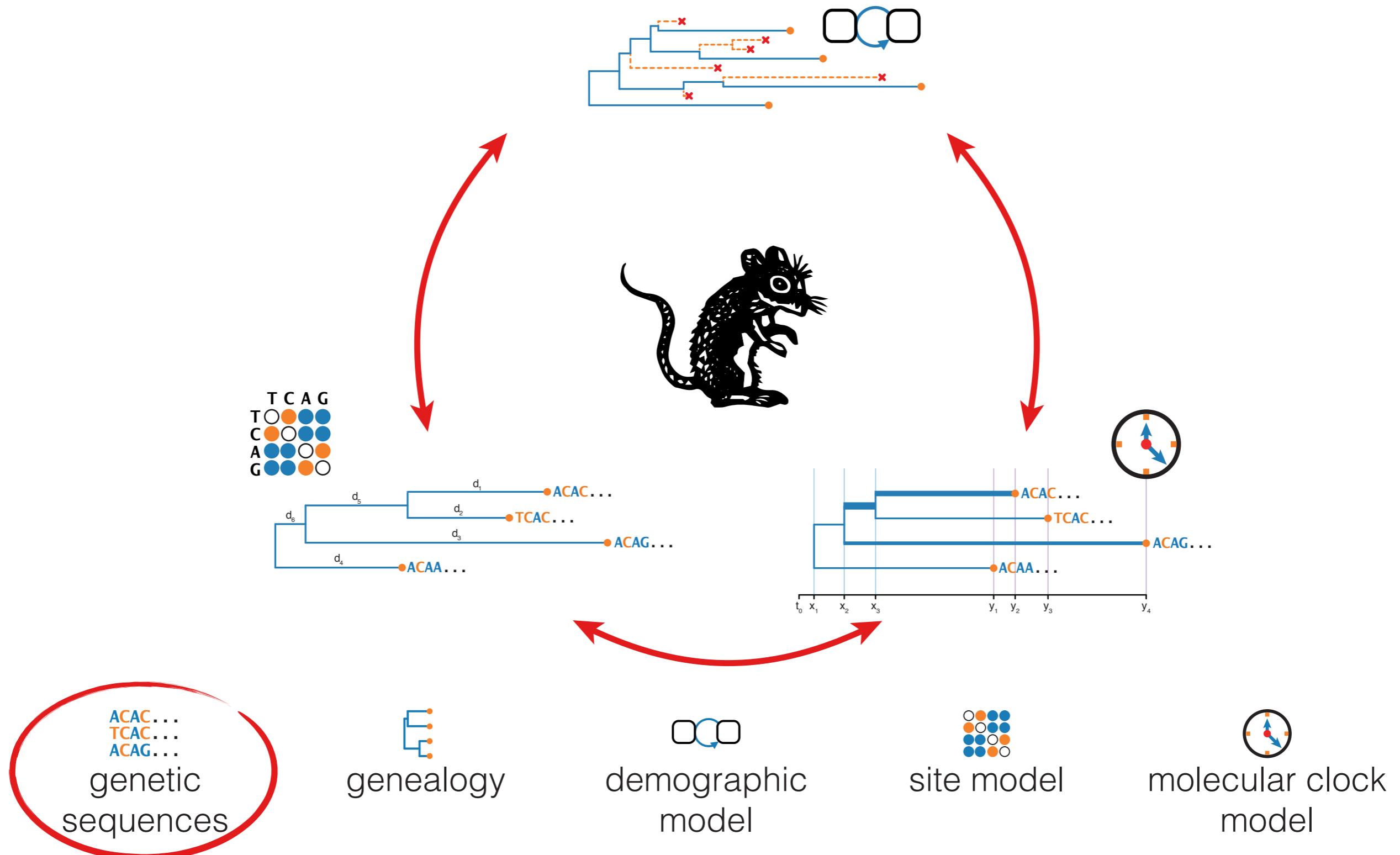
**beast2**      noun

\bēst-tōō\

## Definition of *beast2*

1. **B**ayesian **E**volutionary **A**nalysis by **S**ampling **T**rees **2**
2. a *modular, extensible, cross-platform package for performing Bayesian inference using MCMC with emphasis on phylogenetic analysis of molecular sequences*
3. something formidably difficult to control or deal with

# What goes into a **BEAST** analysis?



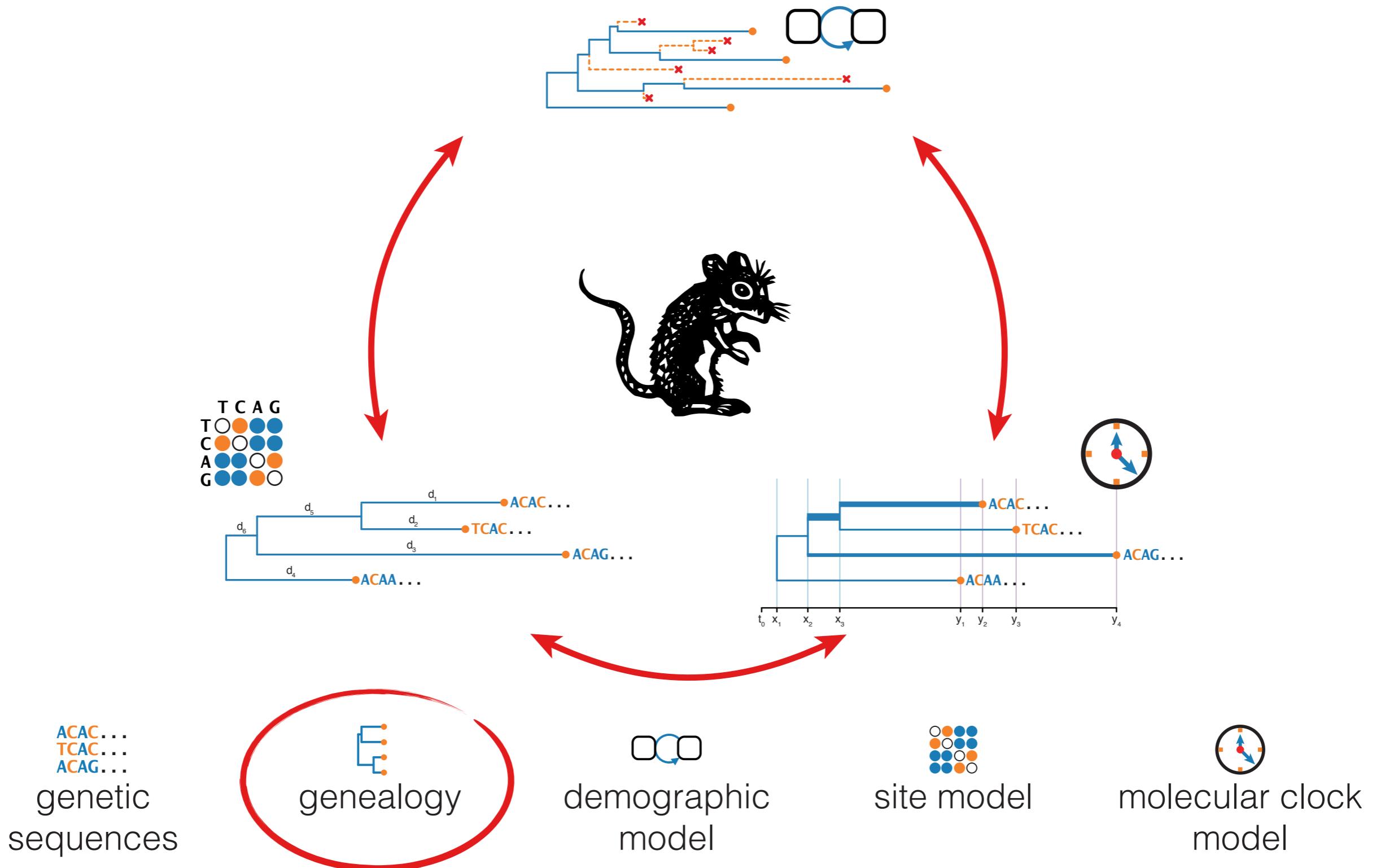
ACAC...  
TCAC...  
ACAG...

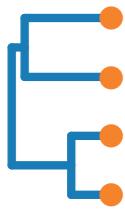
---

## The data

- Samples drawn from a realisation of some stochastic process
- Assume that the data are correct
- Typically one or more alignments of genetic sequencing data (DNA, RNA, amino acids, codons)
- Sampled at one or many time points
- May also contain sampling location or phenotypic trait data

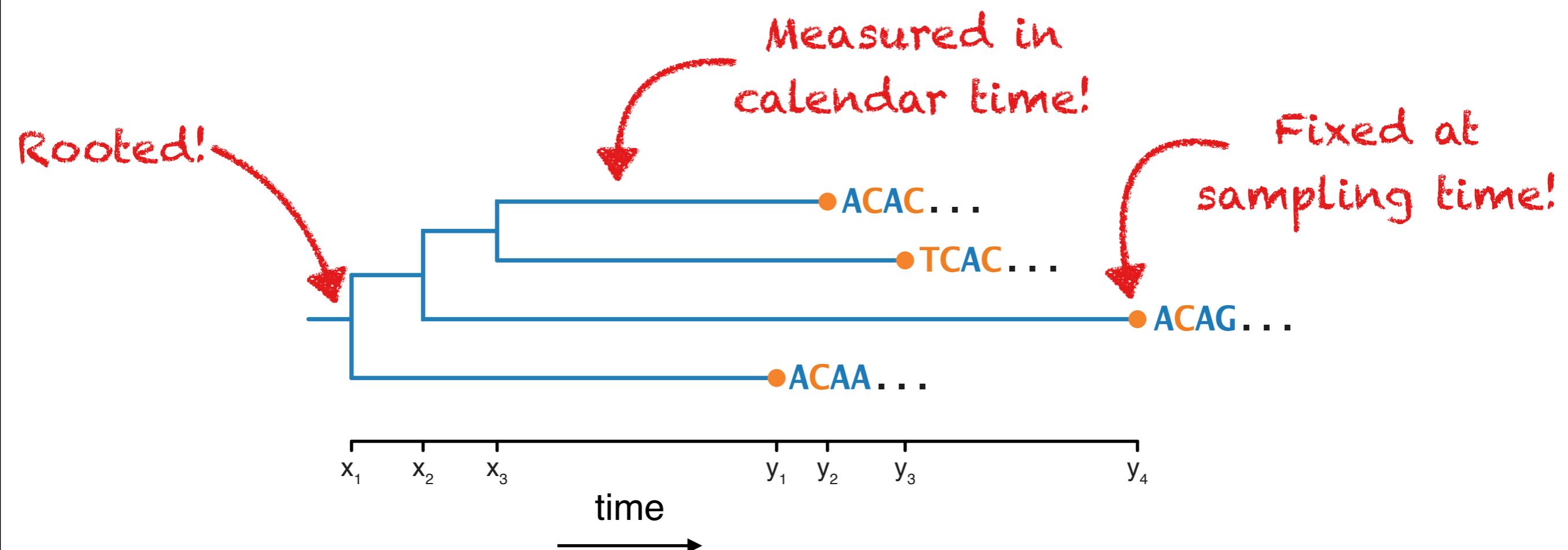
# What goes into a **BEAST** model?



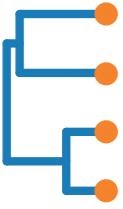


# The genealogy (tree)

The fundamental genealogical structure  
in **BEAST2** is the **rooted time-tree**

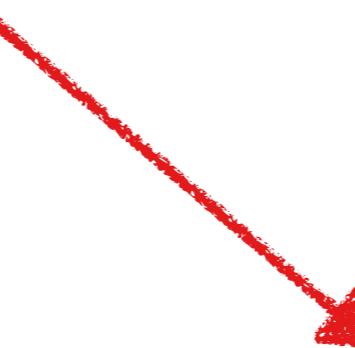
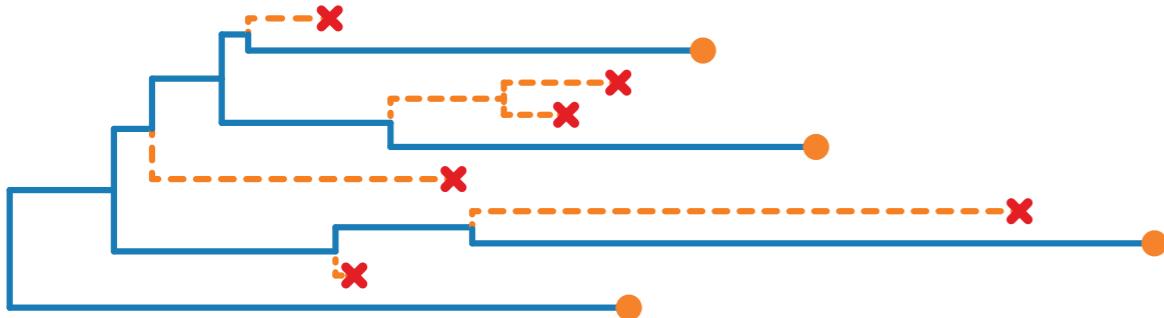


- This tree is a "sampled" or "reconstructed" tree
- Displays ancestral relationships between **sampled sequences** (individuals/taxa)

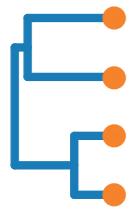


# The genealogy (tree)

**full tree**

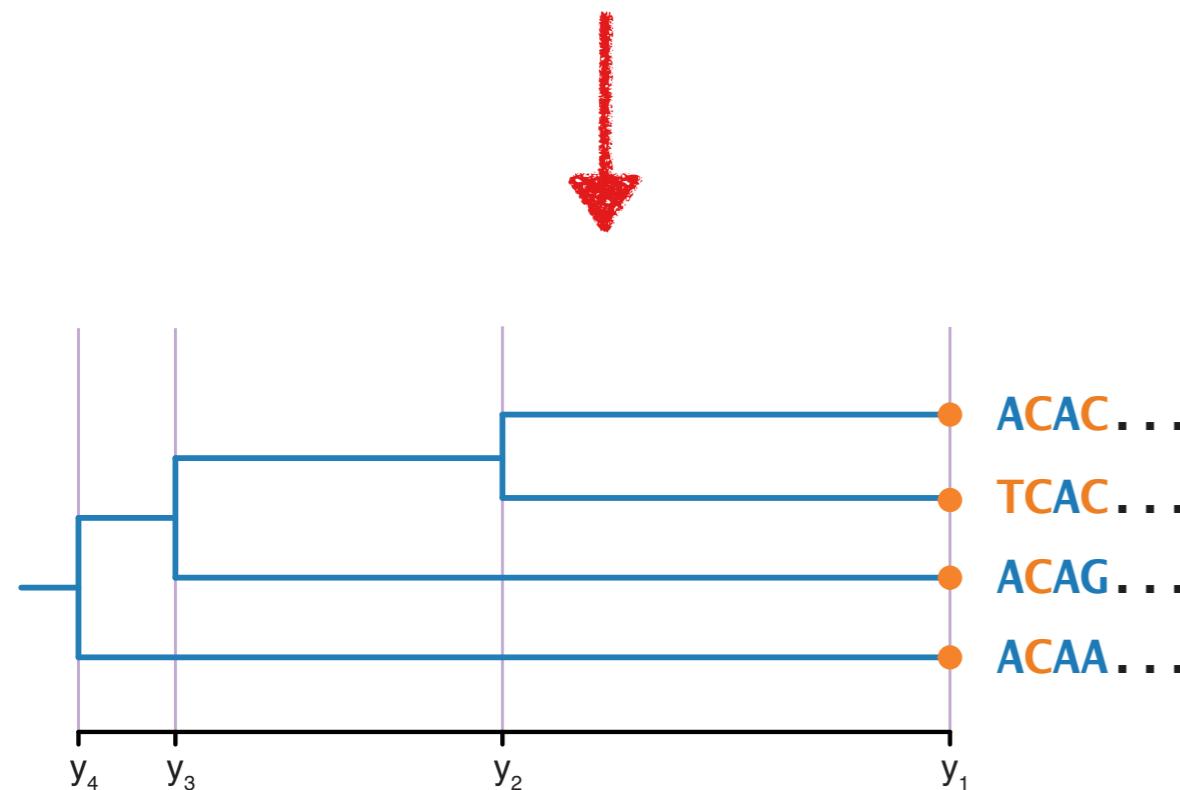


**sampled/reconstructed tree**



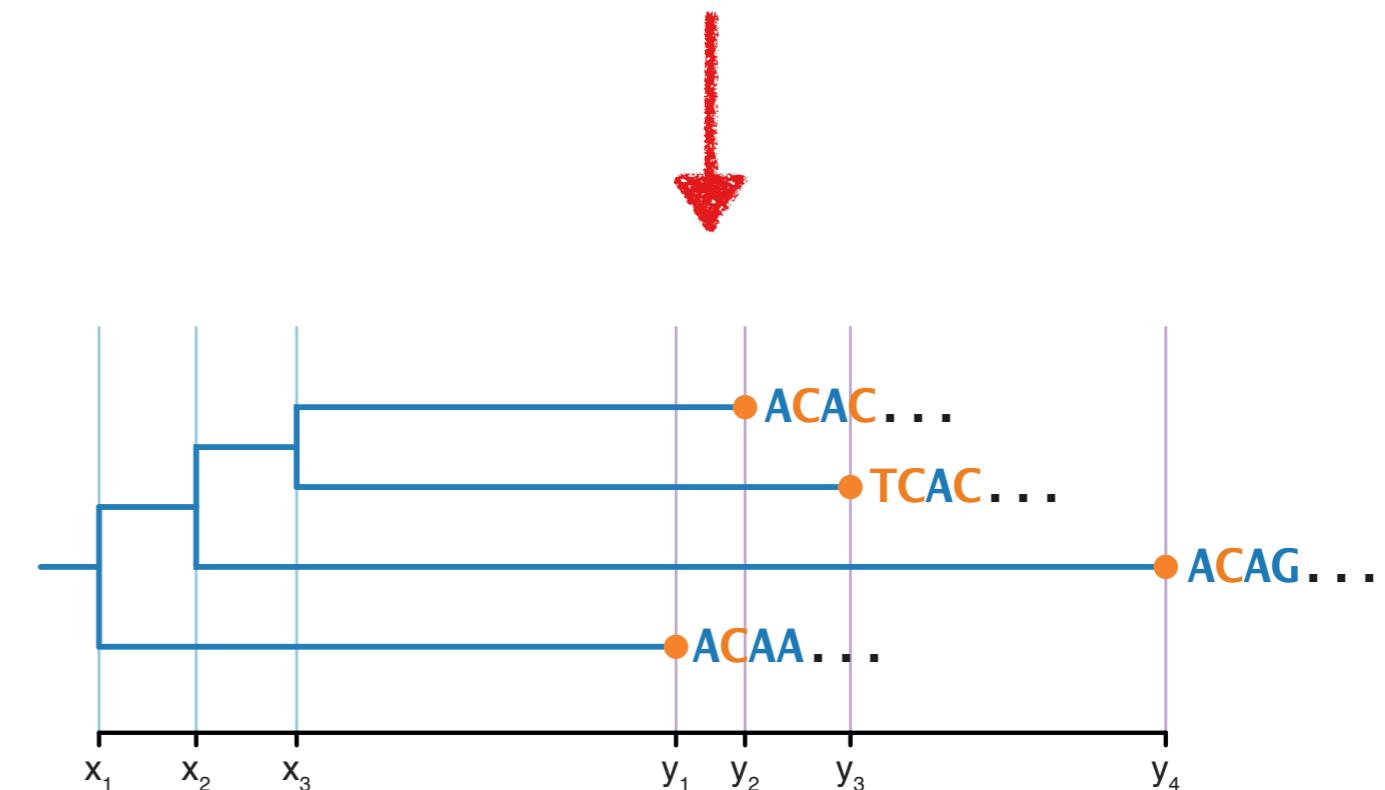
# The genealogy (tree)

Sequences sampled  
at one time point

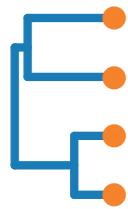


**homochronous tree**

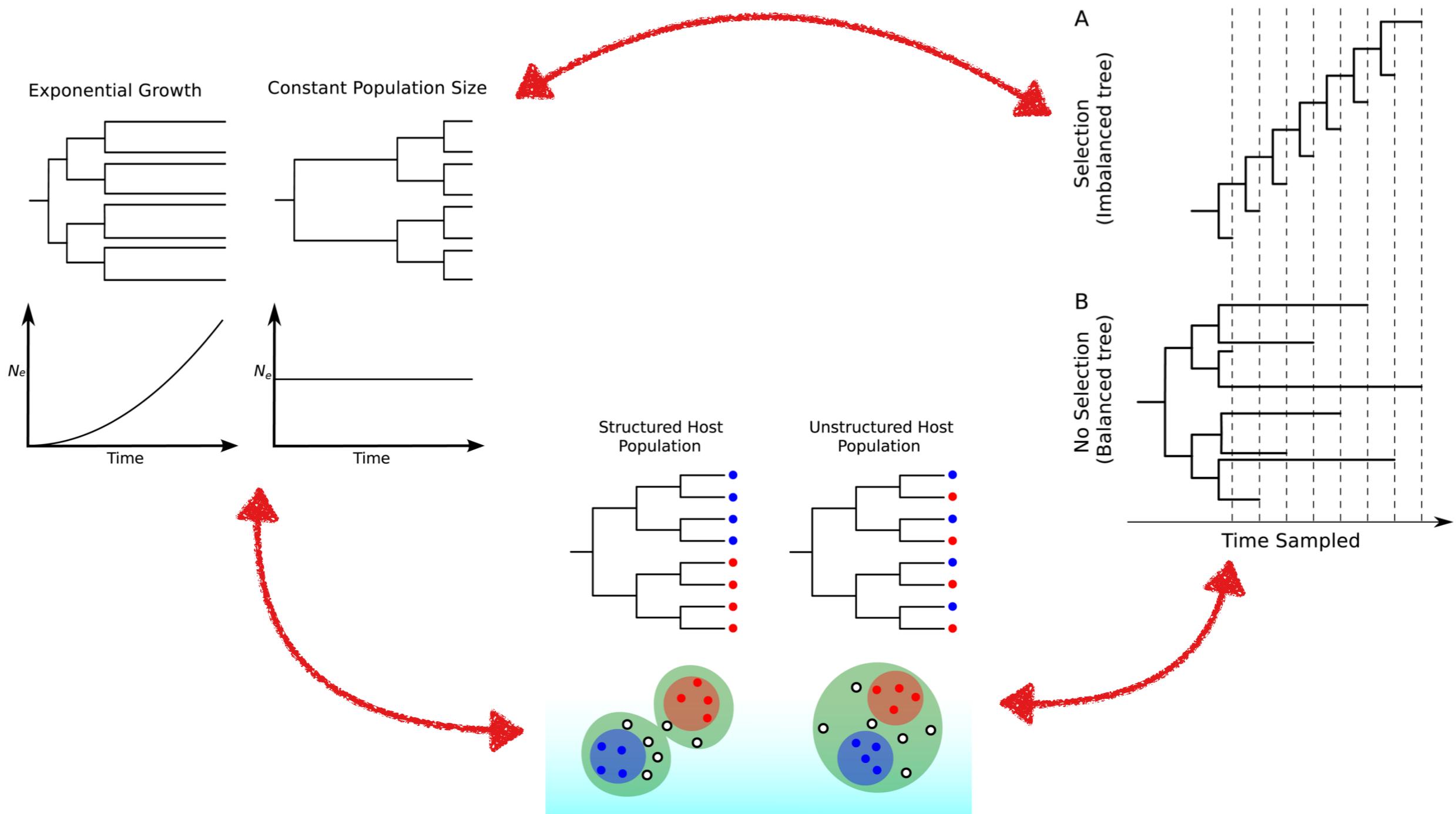
Sequences sampled  
at many time points



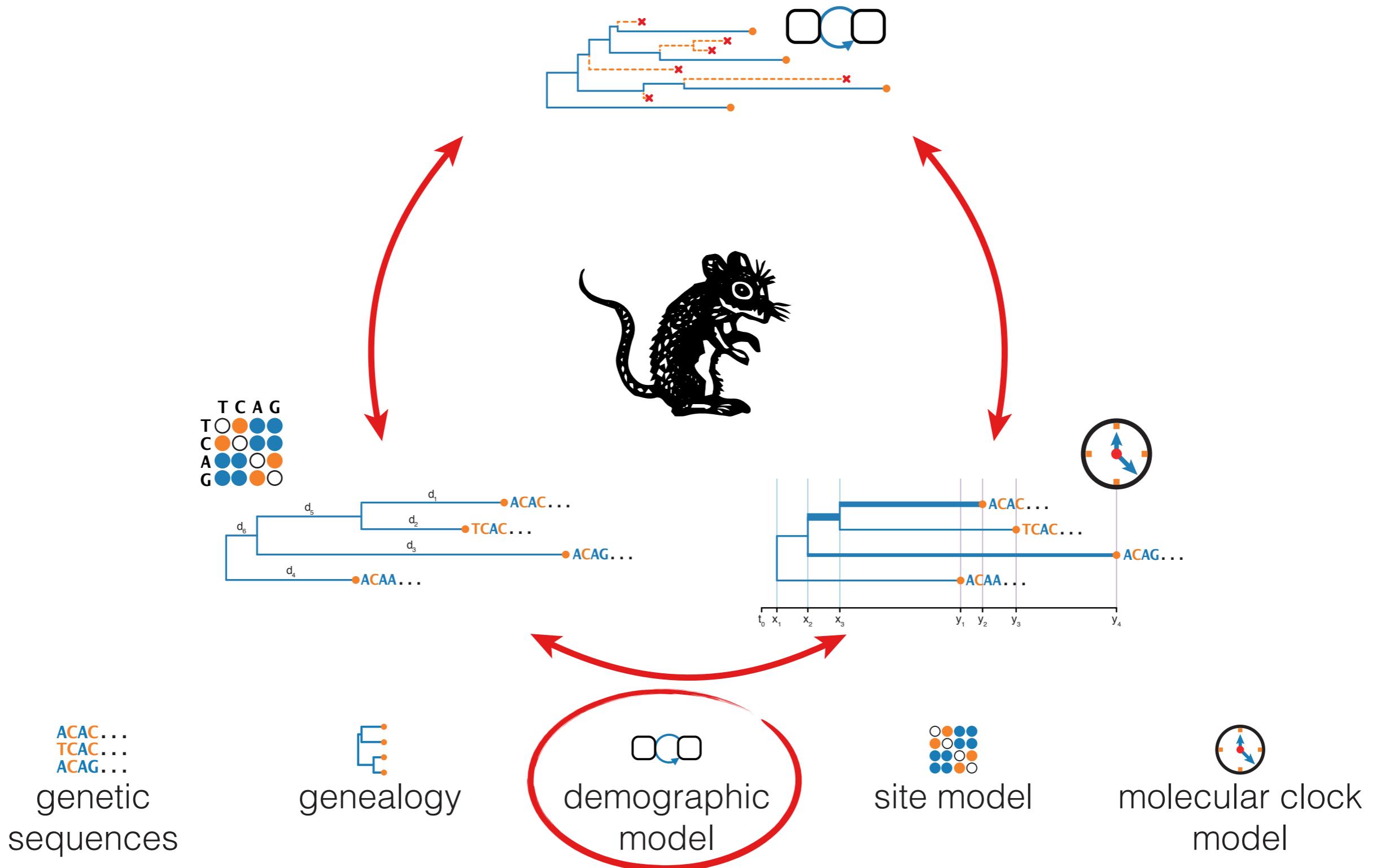
**heterochronous tree**



# Different population dynamics generate trees that look different



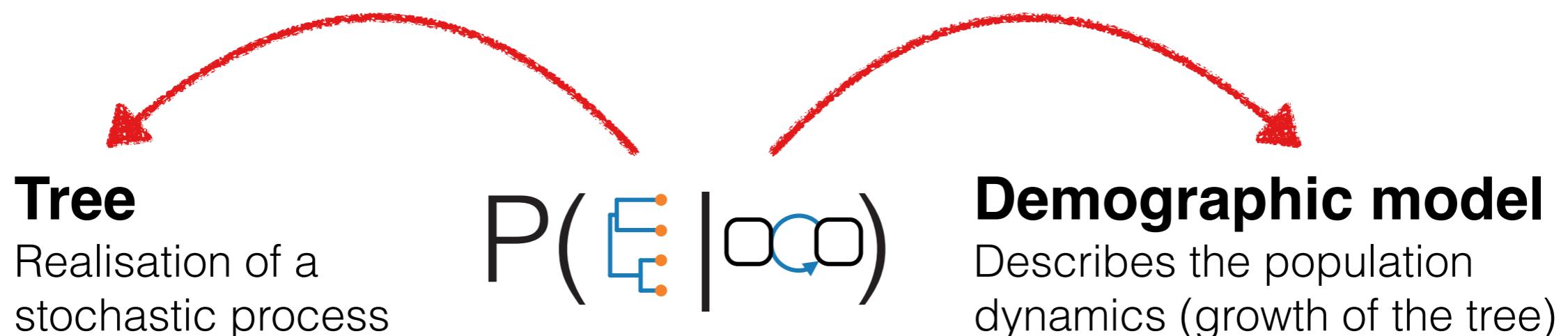
# What goes into a **BEAST** model?





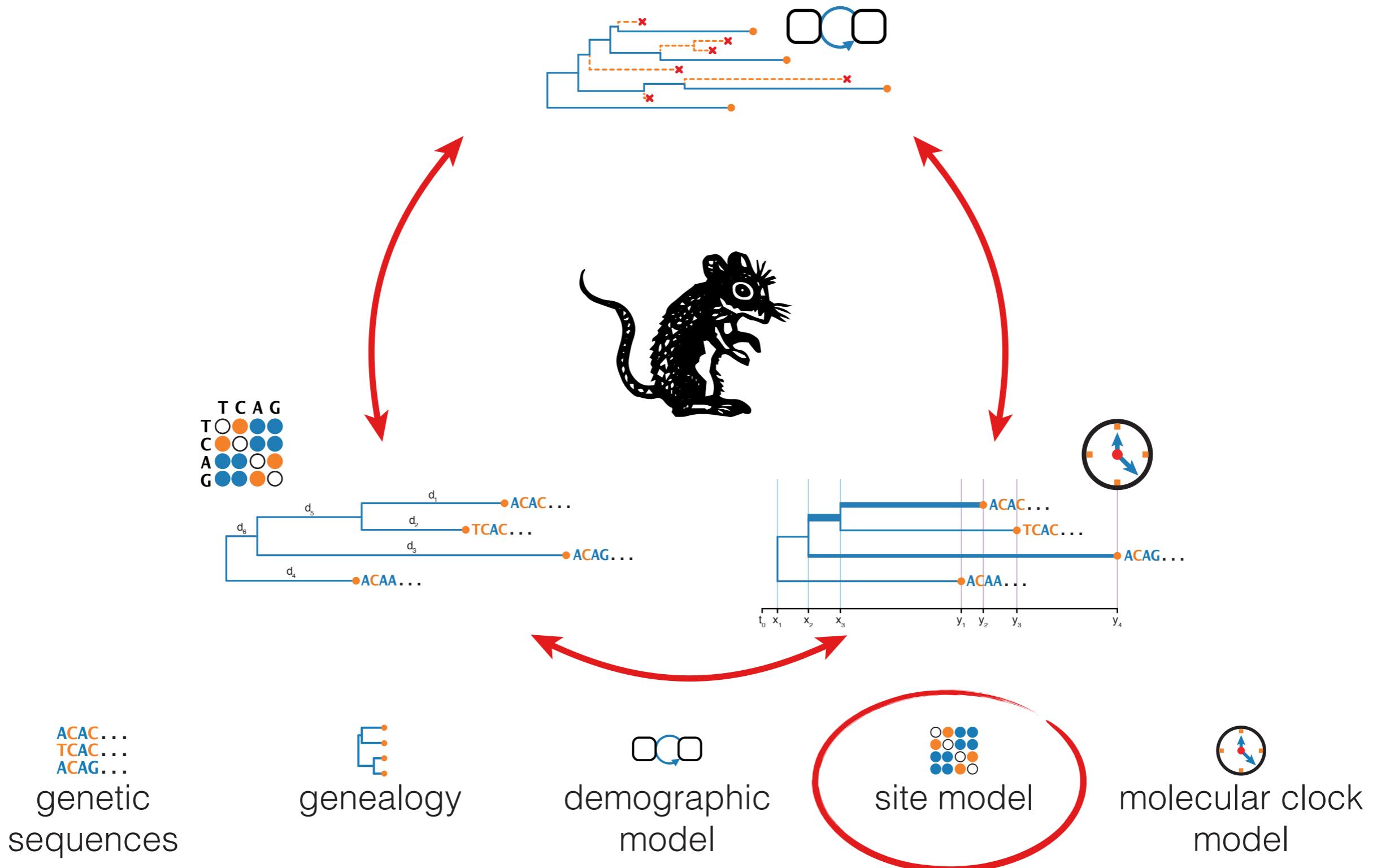
# Demographic model

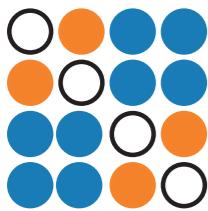
- Describes the population/speciation dynamics
- How does the population demographics / species diversity change over time?



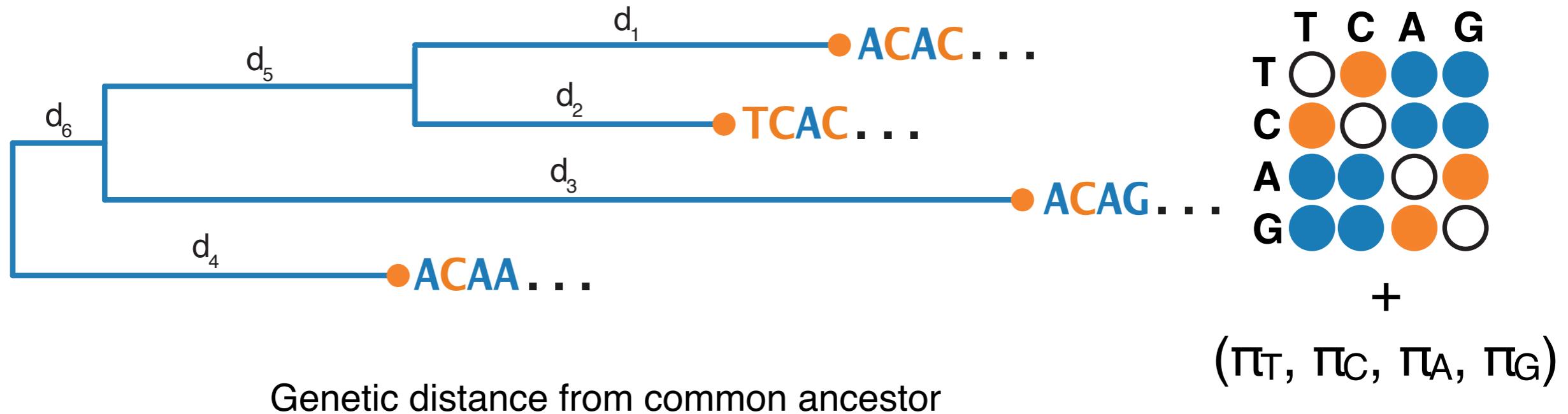
- How likely is the genealogy given a demographic model?
- Sometimes called a **tree prior**
- Usually a **coalescent** or **birth-death** model

# What goes into a **BEAST** model?

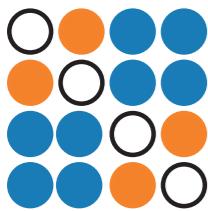




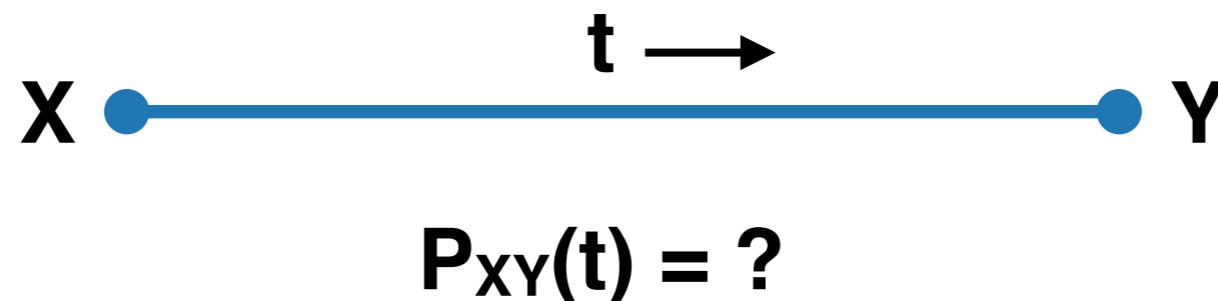
# Site model



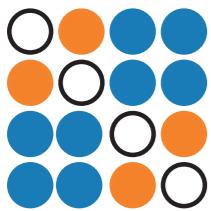
- We observe sequences at the tips, not at internal nodes
- **Substitution model** describes rates of substitution between available characters relative to genetic distance (expected substitutions/site), as well as equilibrium frequencies of characters
- **Site model** describes how the substitution model varies from site-to-site
- **Site model links sequences to the genealogy**
  - using Felsenstein's pruning algorithm we can calculate the likelihood:  $P( \text{ACAC...} | \text{TCAC...} | \text{ACAG...} | \text{ACAA...} )$



# Substitution model



- What is the probability of observing **Y** at the end of the branch?
- Multiple substitutions at the same site means not all substitutions are observed
- Need to account for **all** possible trajectories from **X** to **Y**



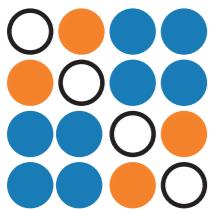
# Substitutions as a Markov process

## Assume:

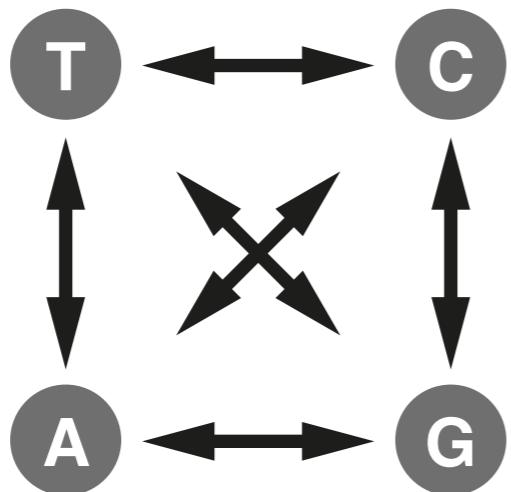
- Every site is evolving independently
- Substitutions at each site is governed by a (usually reversible) Markov process

$$Q = \begin{pmatrix} T & C & A & G \\ T & -(a+b+c) & a & b & c \\ C & d & -(d+e+f) & e & f \\ A & g & h & -(g+h+i) & i \\ G & j & k & l & -(j+k+l) \end{pmatrix}$$

**Process is governed by a rate matrix (Q) which gives the relative rates of substitutions between nucleotides**



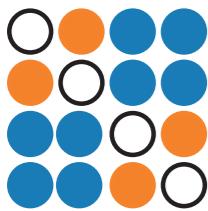
# Jukes-Cantor model (JC69)



$$\begin{matrix} & T & C & A & G \\ T & \cdot & \lambda & \lambda & \lambda \\ C & \lambda & \cdot & \lambda & \lambda \\ A & \lambda & \lambda & \cdot & \lambda \\ G & \lambda & \lambda & \lambda & \cdot \end{matrix}$$

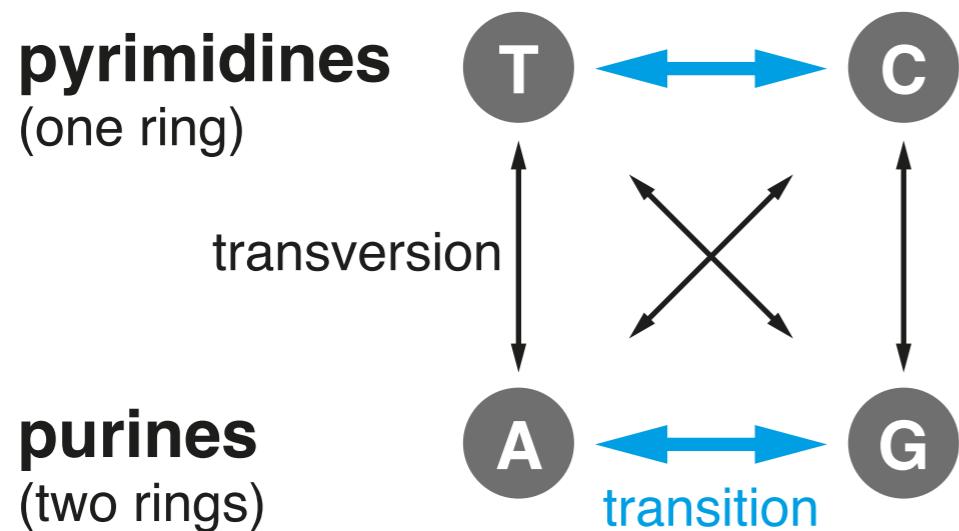
$$\pi_T = \pi_C = \pi_A = \pi_G$$

- Simplest model
- All rates and frequencies are equal!



# Kimura 2-parameter model (K80)

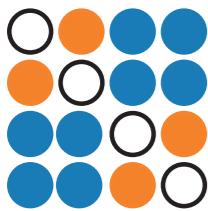
(courtesy of Carsten Magnus)



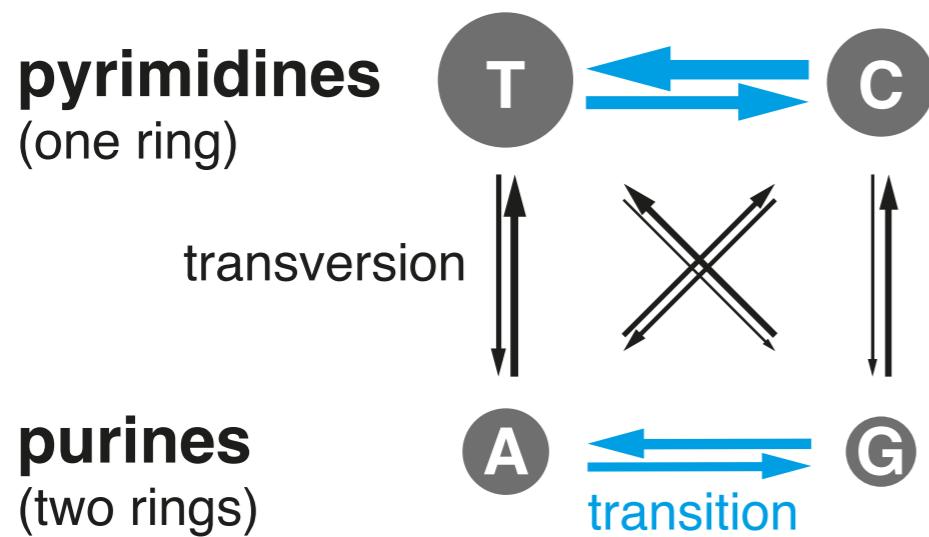
$$\begin{matrix} & T & C & A & G \\ T & \cdot & \alpha & \beta & \beta \\ C & \alpha & \cdot & \beta & \beta \\ A & \beta & \beta & \cdot & \alpha \\ G & \beta & \beta & \alpha & \cdot \end{matrix}$$

$$\pi_T = \pi_C = \pi_A = \pi_G$$

- Accounts for transition/transversion bias
- Still symmetric ( $r_{ij} = r_{ji}$ )
- Equilibrium frequencies still equal  
After a long period of evolution  $p(T) = p(C) = p(A) = p(G) = 0.25$



# HKY-model (HKY85)

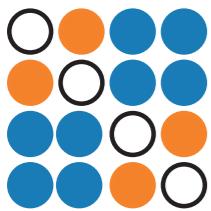


$$(\pi_T, \pi_C, \pi_A, \pi_G)$$

$$\begin{array}{cccc}
 & T & C & A & G \\
 T & \cdot & \alpha\pi_C & \beta\pi_A & \beta\pi_G \\
 C & \alpha\pi_T & \cdot & \beta\pi_A & \beta\pi_G \\
 A & \beta\pi_T & \beta\pi_C & \cdot & \alpha\pi_G \\
 G & \beta\pi_T & \beta\pi_C & \alpha\pi_A & \cdot
 \end{array}$$

$$= \begin{pmatrix} \cdot & \alpha & \beta & \beta \\ \alpha & \cdot & \beta & \beta \\ \beta & \beta & \cdot & \alpha \\ \beta & \beta & \alpha & \cdot \end{pmatrix} \cdot \begin{pmatrix} \pi_T & 0 & 0 & 0 \\ 0 & \pi_C & 0 & 0 \\ 0 & 0 & \pi_A & 0 \\ 0 & 0 & 0 & \pi_G \end{pmatrix}$$

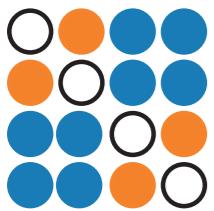
- Accounts for transition/transversion bias
- Accounts for unequal equilibrium frequencies
- Not symmetric anymore ( $r_{ij} \neq r_{ji}$ )
- Still time-reversible ( $\pi_i q_{ij} = \pi_j q_{ji}$ )



# General time-reversible model (GTR/REV) (courtesy of Carsten Magnus)

$$\begin{array}{cccc} T & C & A & G \\ \hline T & \cdot & a\pi_C & b\pi_A & c\pi_G \\ C & a\pi_T & \cdot & d\pi_A & e\pi_G \\ A & b\pi_T & d\pi_C & \cdot & f\pi_G \\ G & c\pi_T & e\pi_C & f\pi_A & \cdot \end{array} = \begin{pmatrix} \cdot & a & b & c \\ a & \cdot & d & e \\ b & d & \cdot & f \\ c & e & f & \cdot \end{pmatrix} \cdot \begin{pmatrix} \pi_T & 0 & 0 & 0 \\ 0 & \pi_C & 0 & 0 \\ 0 & 0 & \pi_A & 0 \\ 0 & 0 & 0 & \pi_G \end{pmatrix}$$

- Most general time-reversible model
- More flexible models are possible, but mathematically inconvenient



# Transition probability matrix

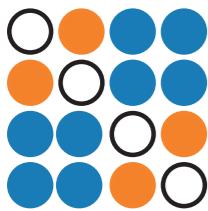
$$\mathbf{P}(t) = e^{\mathbf{Q}t} \quad \mathbf{P}(t) = \begin{pmatrix} T & C & A & G \\ T & p_{tt}(t) & p_{tc}(t) & p_{ta}(t) & p_{tg}(t) \\ C & p_{ct}(t) & p_{cc}(t) & p_{ca}(t) & p_{cg}(t) \\ A & p_{at}(t) & p_{ac}(t) & p_{aa}(t) & p_{ag}(t) \\ G & p_{gt}(t) & p_{gc}(t) & p_{ga}(t) & p_{gg}(t) \end{pmatrix}$$

- Transition probabilities take into account every possible evolutionary trajectory (Chapman-Kolmogorov theorem)

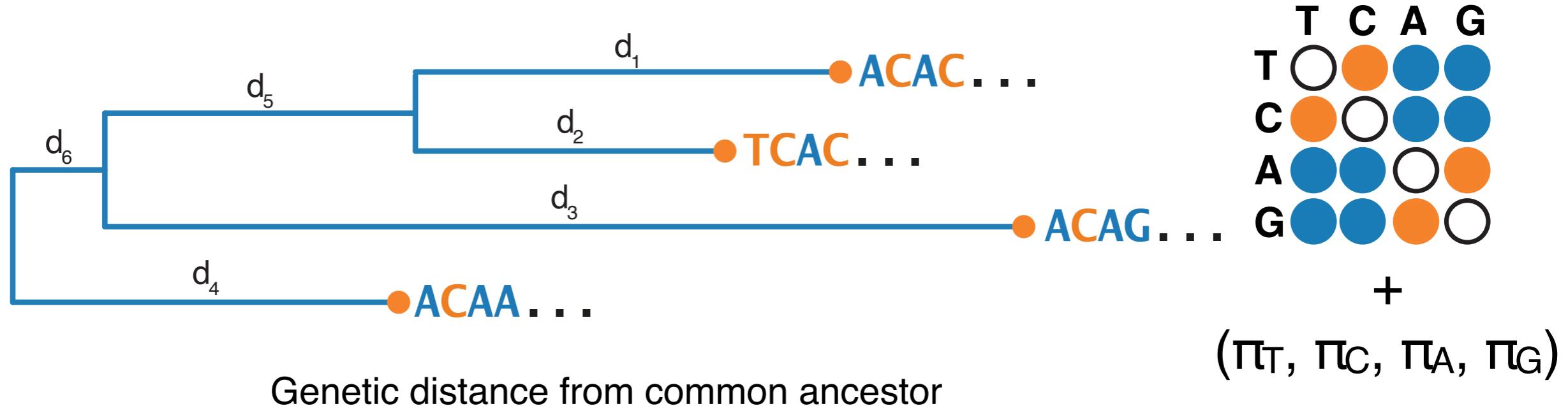


$$P_{XY}(t) = ?$$

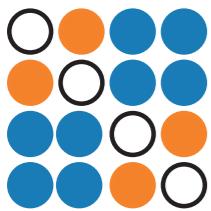
- $\mathbf{Q}$  only gives the **relative substitution rates**  
⇒ **distance is measured in expected number of substitutions per site**



# Site model



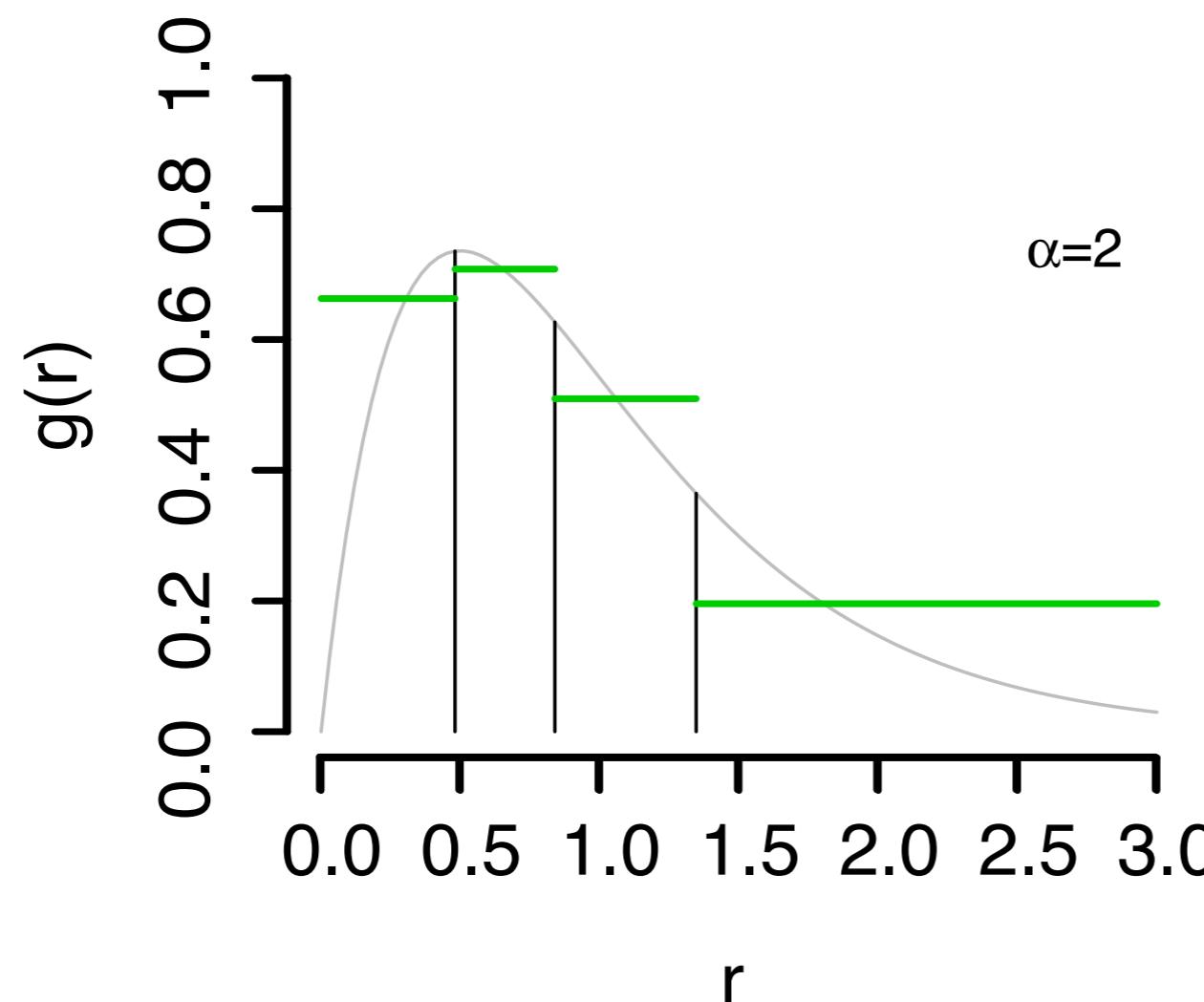
- Describes how the substitution model varies from site-to-site
- Assume every site is evolving independently
- Account for rate heterogeneity between sites:
  - Proportion of invariant sites
  - Gamma rate heterogeneity
  - Multi-locus models

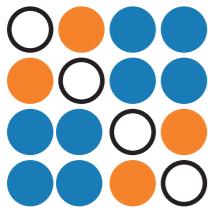


# Gamma rate heterogeneity

(courtesy of Carsten Magnus)

- Not all sites evolve at the same rate
- Assume rate heterogeneity among sites is  $\Gamma$ -distributed
- Discretise  $\Gamma$ -distribution to  $n$  discrete rate categories for computational reasons



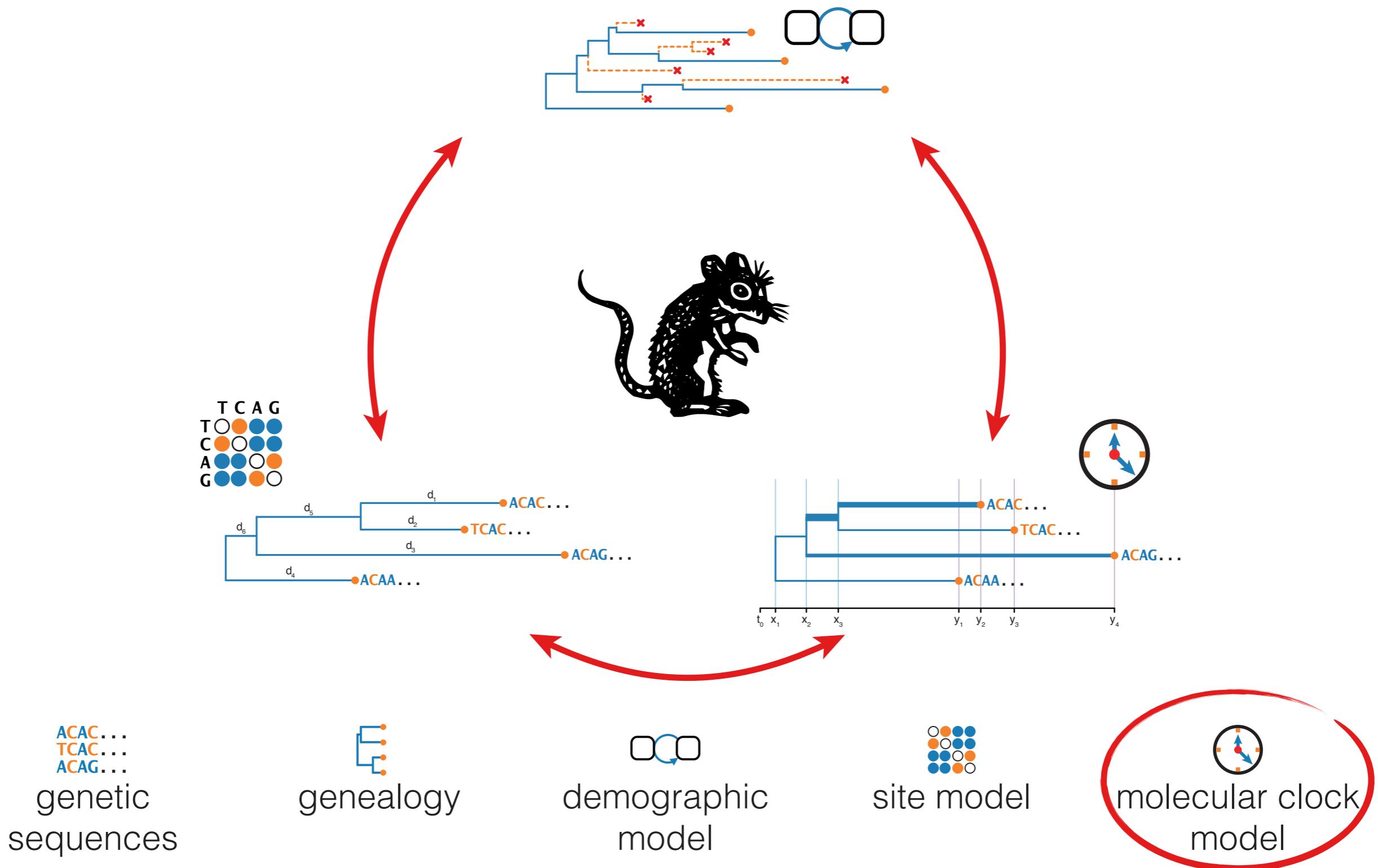


# Multi-locus models

---

- $\Gamma$ -distributed rate variation is not always flexible enough to model differences between different loci
- Use a separate substitution model for each locus
- Can also use separate models for different codon positions
- Loci could be different genes, different codon positions, etc.

# What goes into a **BEAST** model?



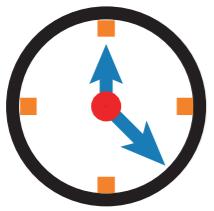
genetic  
sequences

genealogy

demographic  
model

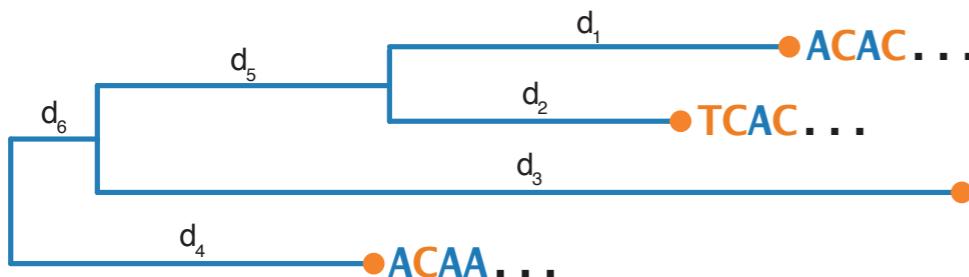
site model

molecular clock  
model



# Molecular clock model

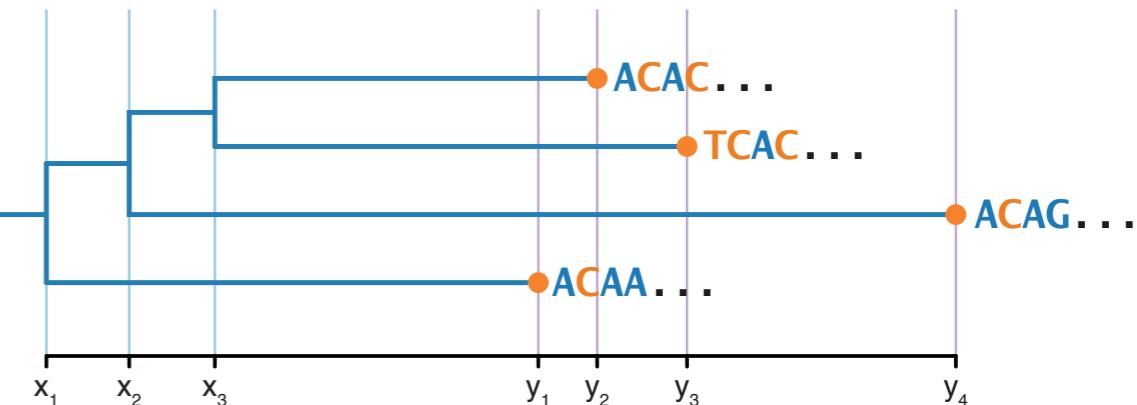
**genetic distance tree**  
(subst/site)



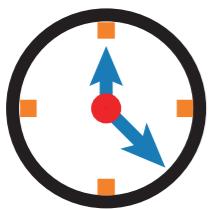
**clock rate**  
(subst/site/year)

$$= \mu \times$$

**time tree**  
(years)

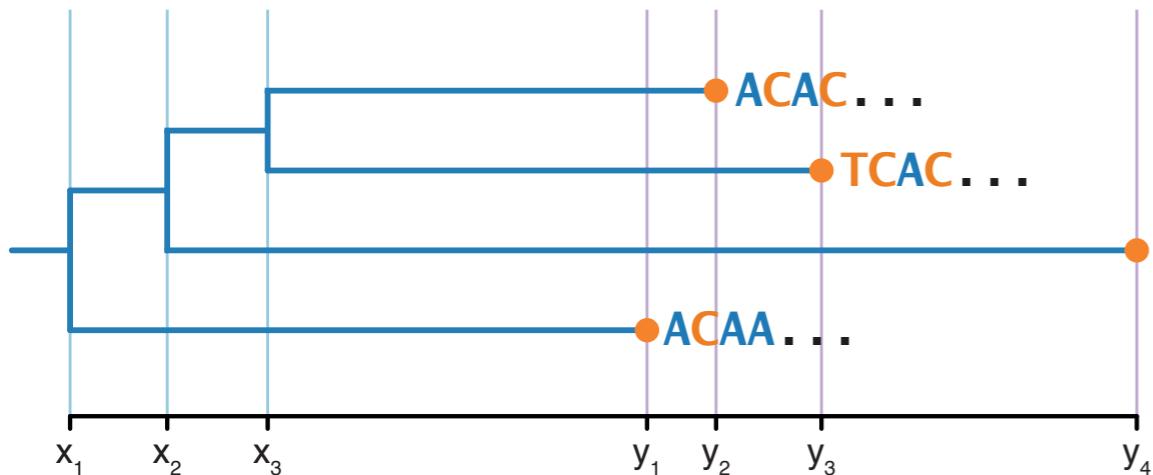


- Determines how quickly sequences are evolving along the tree
- **Genetic distance = Rate x Time**

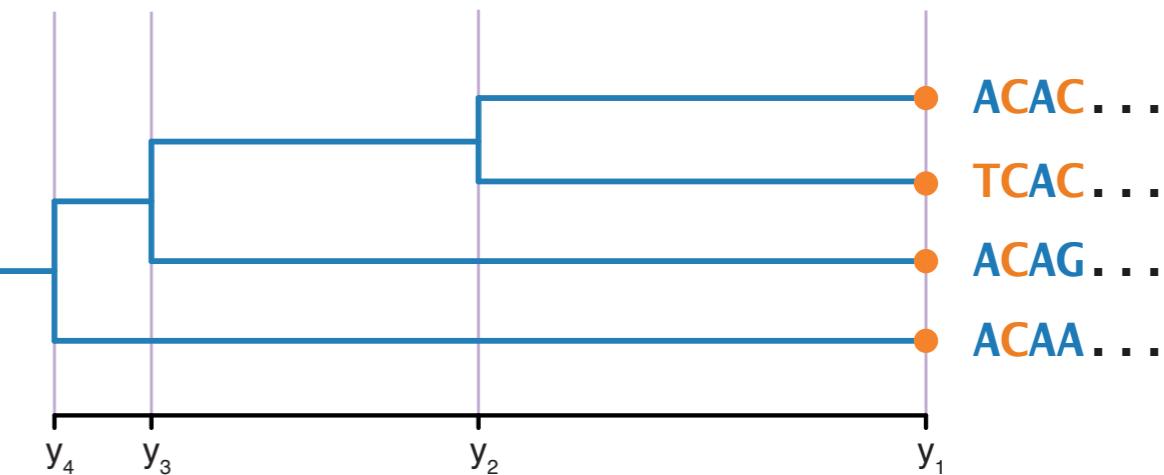


# Molecular clock model

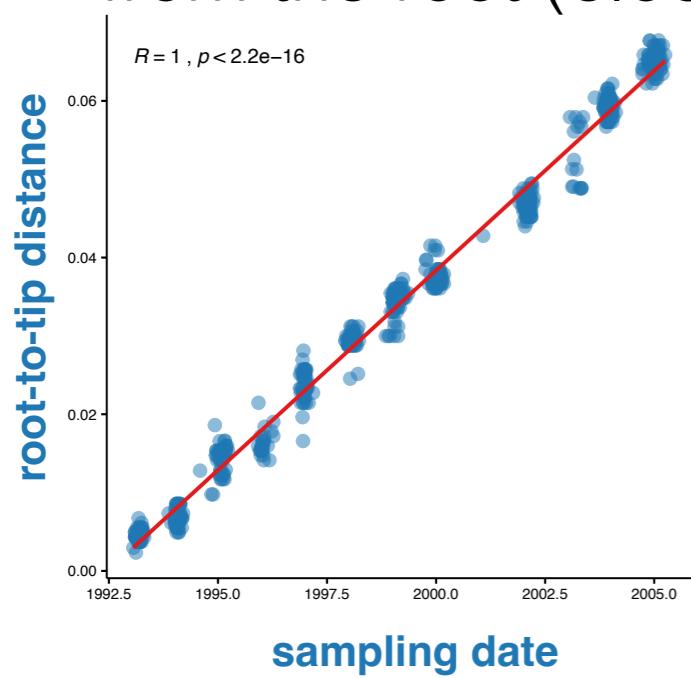
## heterochronous tree



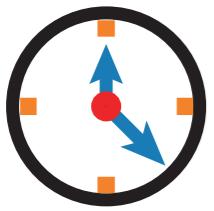
## homochronous tree



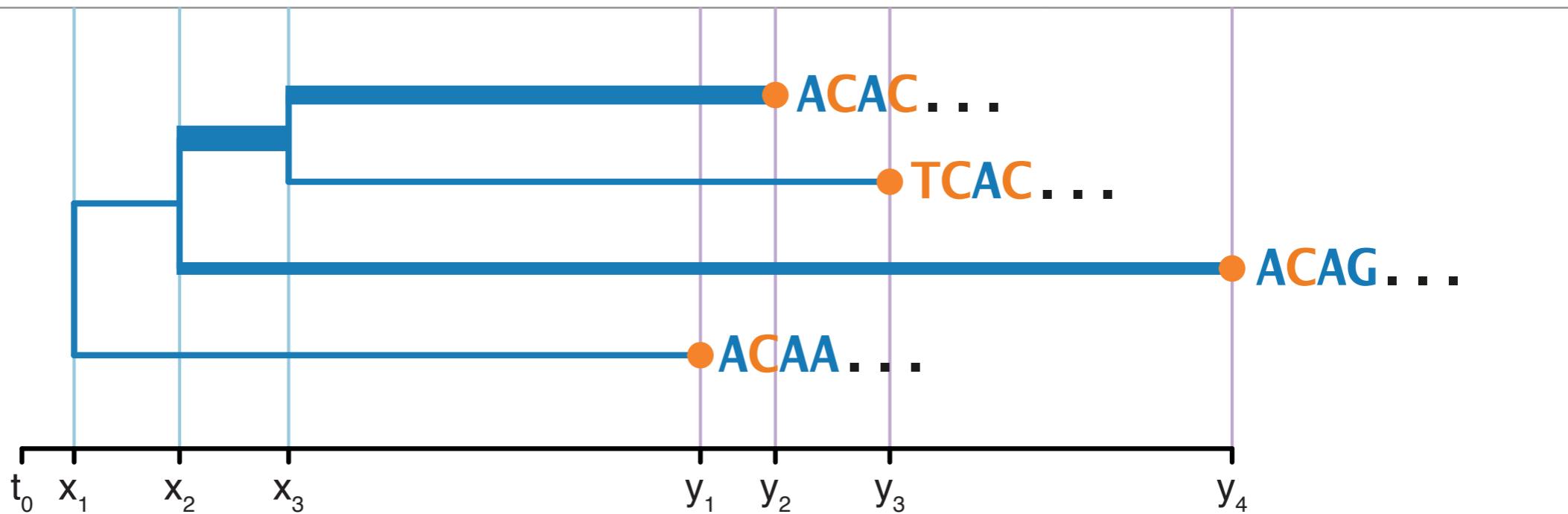
- Correlation between sampling time and genetic distance from the root (clock signal)



- Need external information to calibrate the clock
- Fix clock rate or calibrate some internal nodes with external information



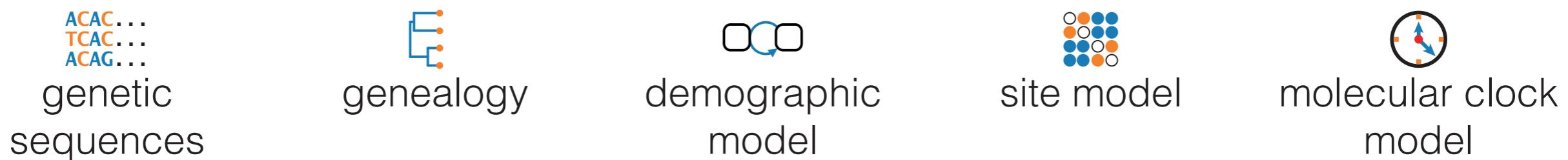
# Relaxed and local clocks



Some branches may have different clock rates!

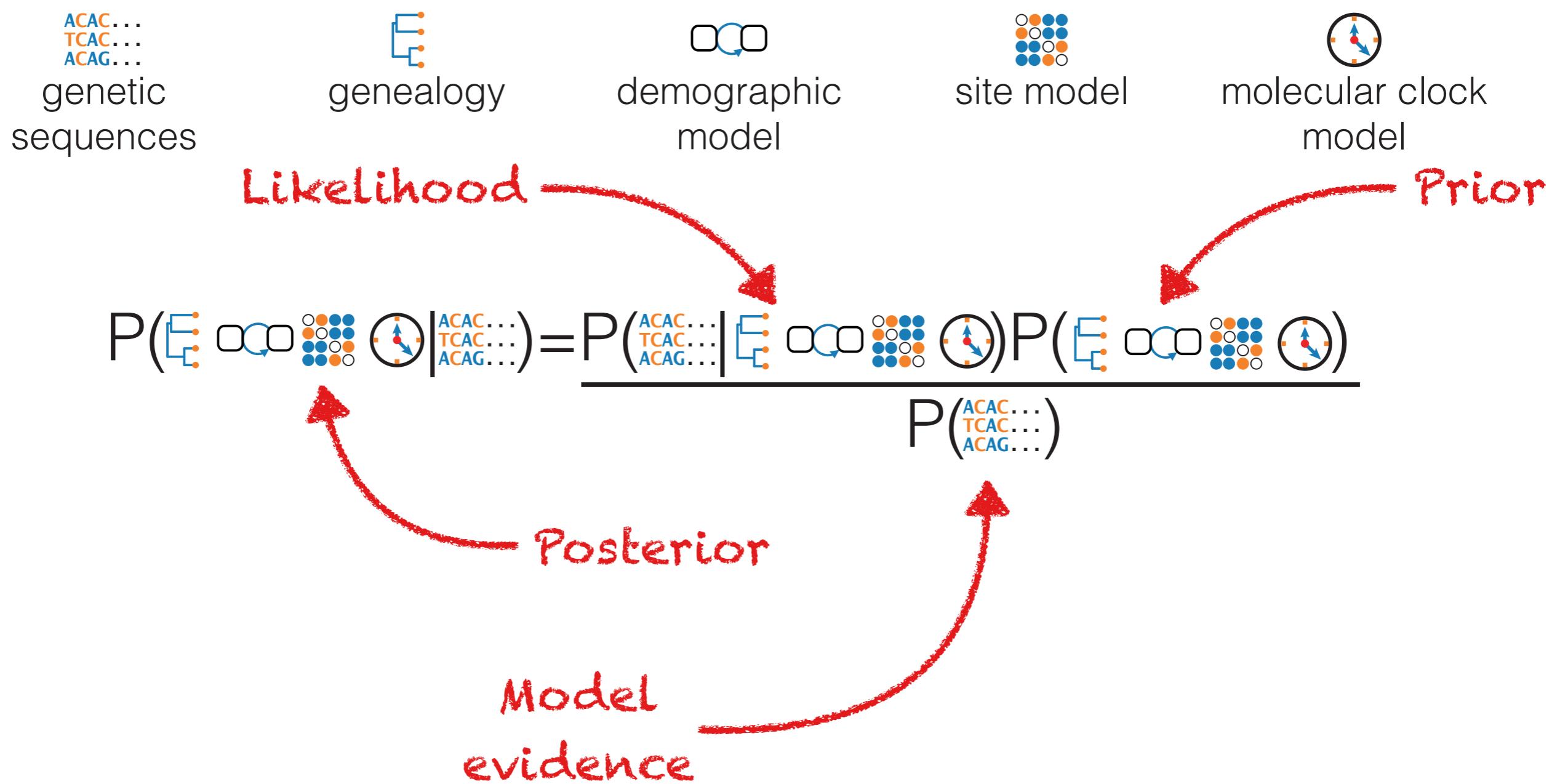
# Putting it all together

---



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

# Putting it all together



# Putting it all together

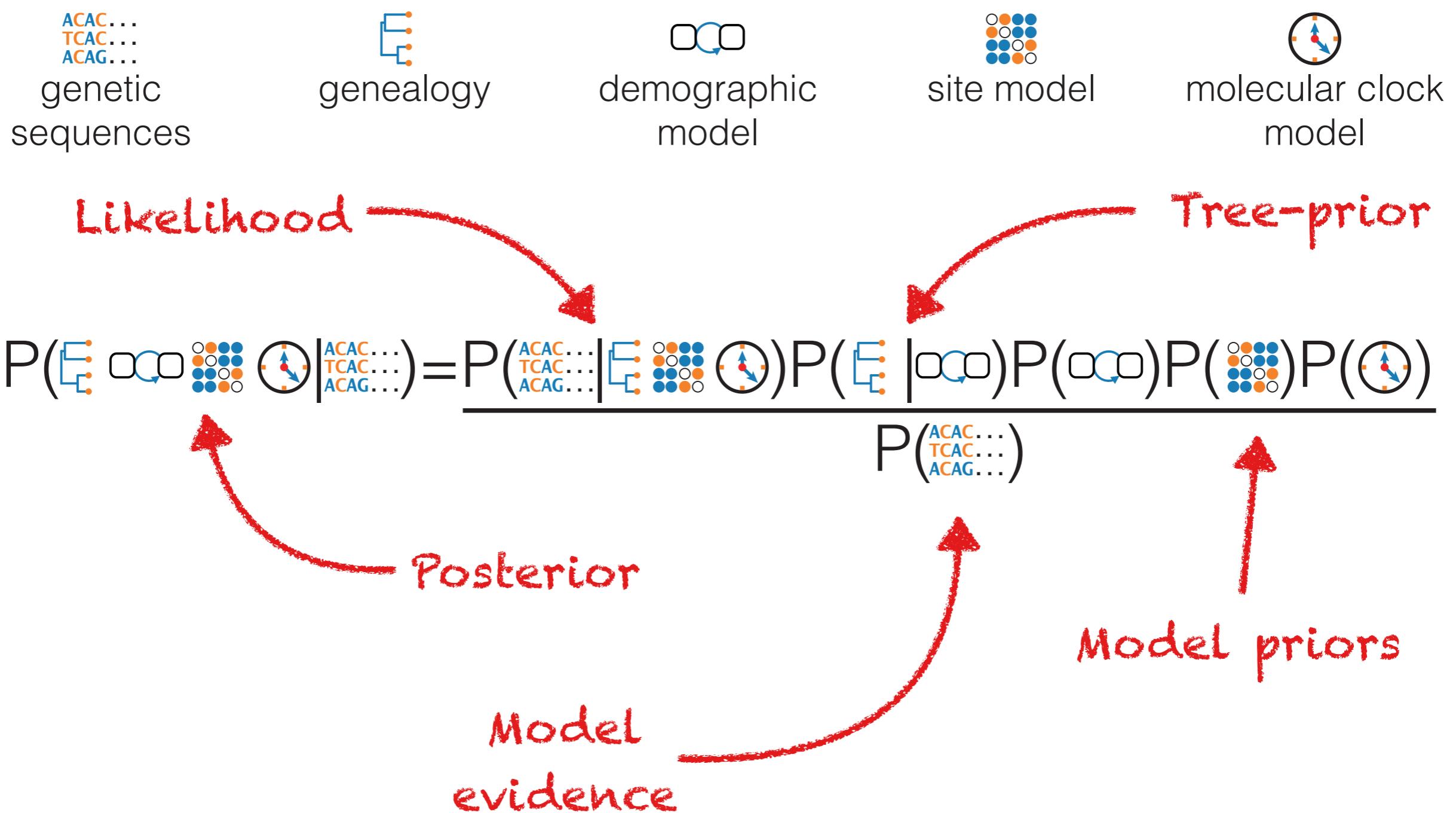


$$P(\text{E} \mid \text{ACAC...}, \text{TCAC...}, \text{ACAG...}) = \frac{P(\text{ACAC...} \mid \text{E}) P(\text{TCAC...} \mid \text{E}) P(\text{ACAG...} \mid \text{E})}{P(\text{ACAC...})}$$

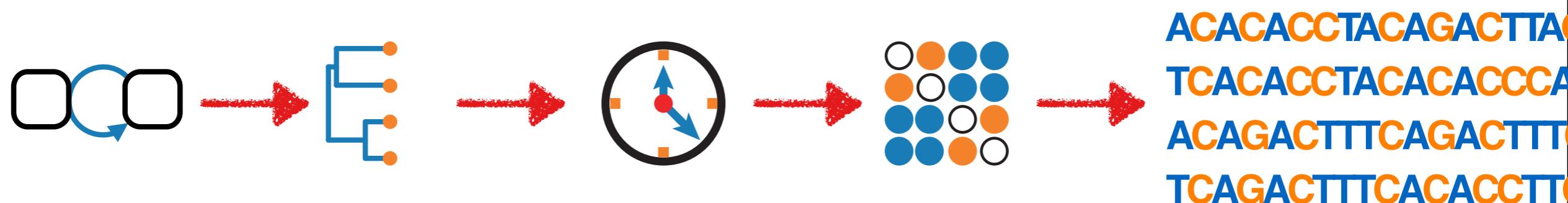
**Assume independence**

$$P(\text{E} \mid \text{ACAC...}, \text{TCAC...}, \text{ACAG...}) = P(\text{E} \mid \text{ACAC...}) P(\text{E} \mid \text{TCAC...}) P(\text{E} \mid \text{ACAG...})$$

# Posterior distribution in BEAST2



# Neutrality assumption



**Assumed tree-generating process  
independent of sequence data**

$$P(E \circ \square \circ \bullet \bullet \bullet | \text{ACAC} \dots) = \frac{P(\text{ACAC} \dots | E \circ \square \circ \bullet \bullet \bullet) P(E | \circ \square) P(\circ \square) P(\bullet \bullet \bullet) P(\bullet \bullet \bullet)}{P(\text{ACAC} \dots)}$$

# How can we find the posterior?

- We want to calculate the posterior distribution

$$P(E \cap O \cap D \cap C | ACAC, TCAC, ACAG, \dots) =$$


- But we cannot easily calculate the marginal likelihood (model evidence)

$$P(A C A C \dots) \rightarrow ?$$

→ use **MCMC!** (Markov-chain Monte Carlo)

- MCMC is a stochastic algorithm that performs a random walk on the posterior, preferentially sampling high-density areas

# MCMC

(Markov-chain Monte Carlo)

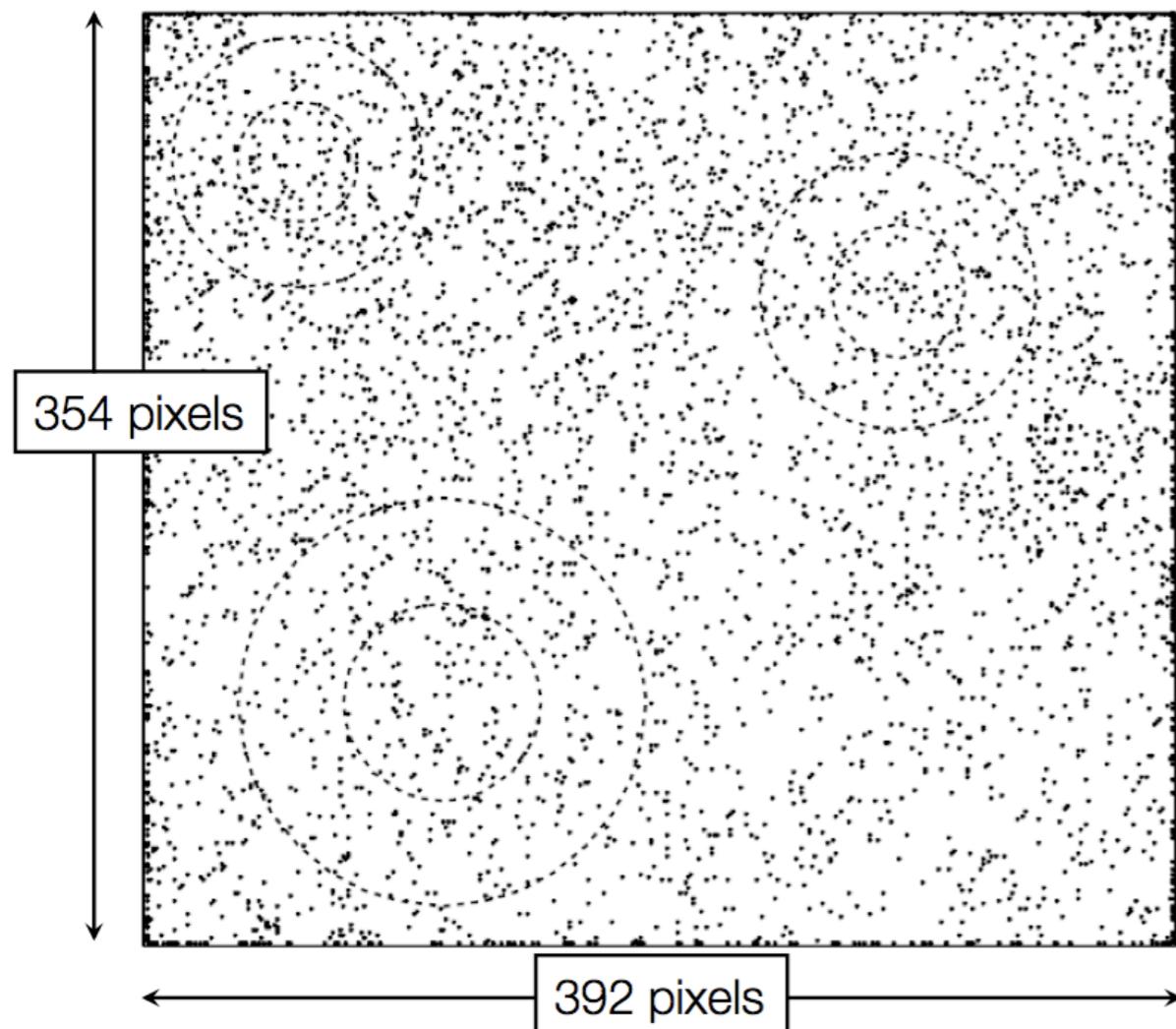
---

- MCMC draws samples from the posterior
  - output is a list of values that can approximate the posterior
- Only need to compare which posterior density is higher
  - So we only need the ratio of posteriors  
(marginal likelihoods cancel out!)

$$\frac{P(\theta_1 | D)}{P(\theta_2 | D)} = \frac{\frac{P(D | \theta_1)P(\theta_1)}{P(D)}}{\frac{P(D | \theta_2)P(\theta_2)}{P(D)}}$$

# Pure random walk (courtesy of Paul Lewis)

---



## **Random walk**

- Random direction
- Gamma distributed step size
- Reflection at edges

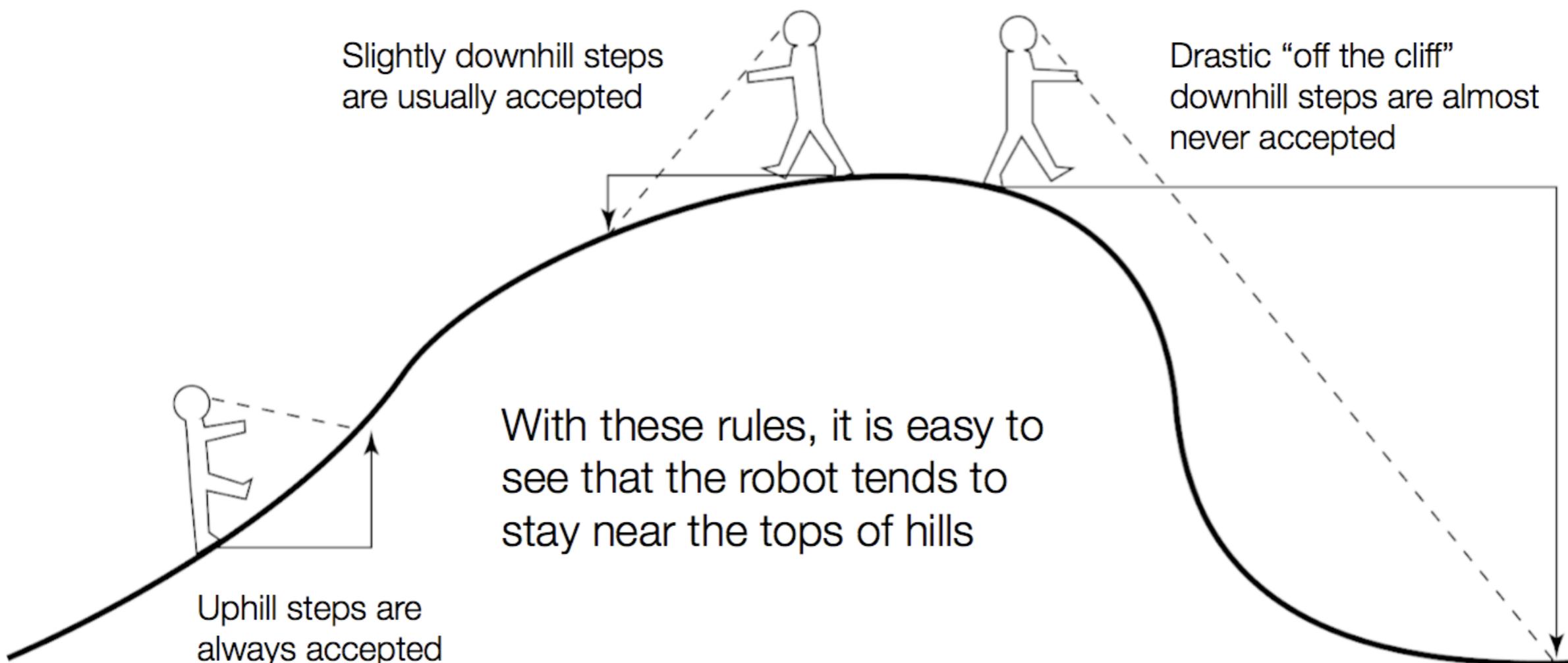
## **Target distribution**

- Equal mixture of 3 bivariate normal hills
- Inner contours: 50%
- Outer contours: 95%

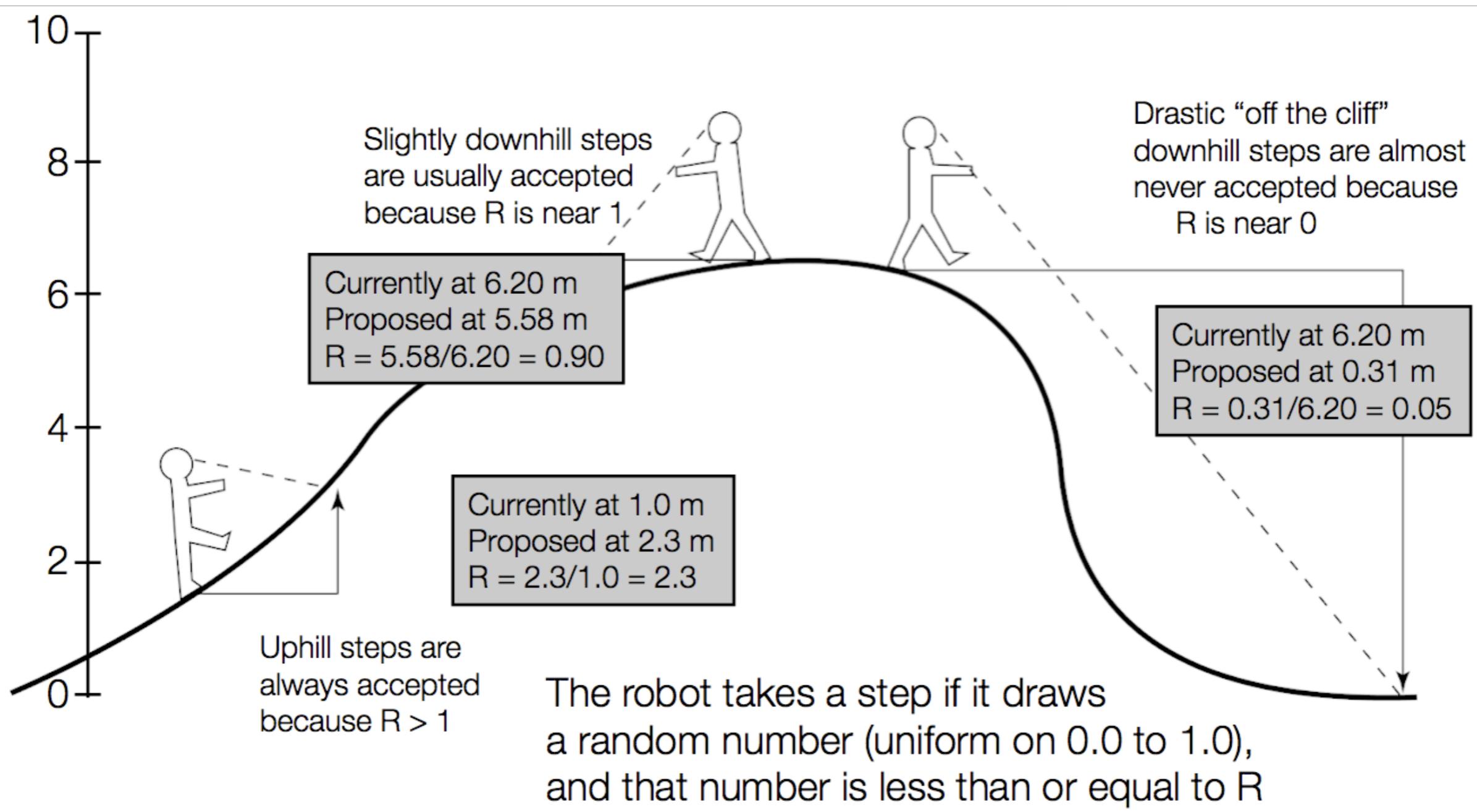
5000 steps by the random walk - not informative at all!

# MCMC robot (courtesy of Paul Lewis)

---



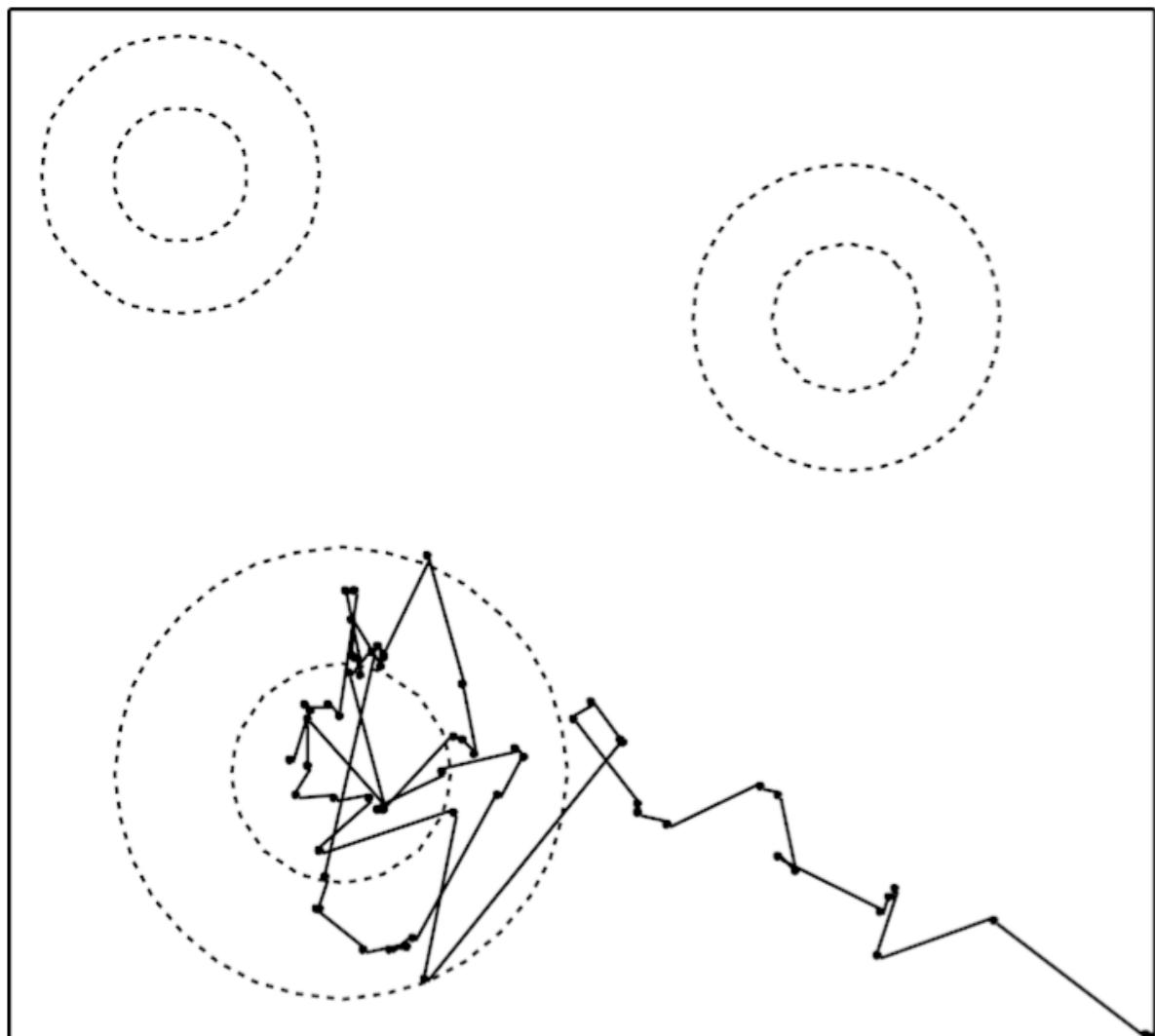
# MCMC robot (courtesy of Paul Lewis)



( $R$  is the ratio between the posterior densities)

# Burn in (courtesy of Paul Lewis)

---

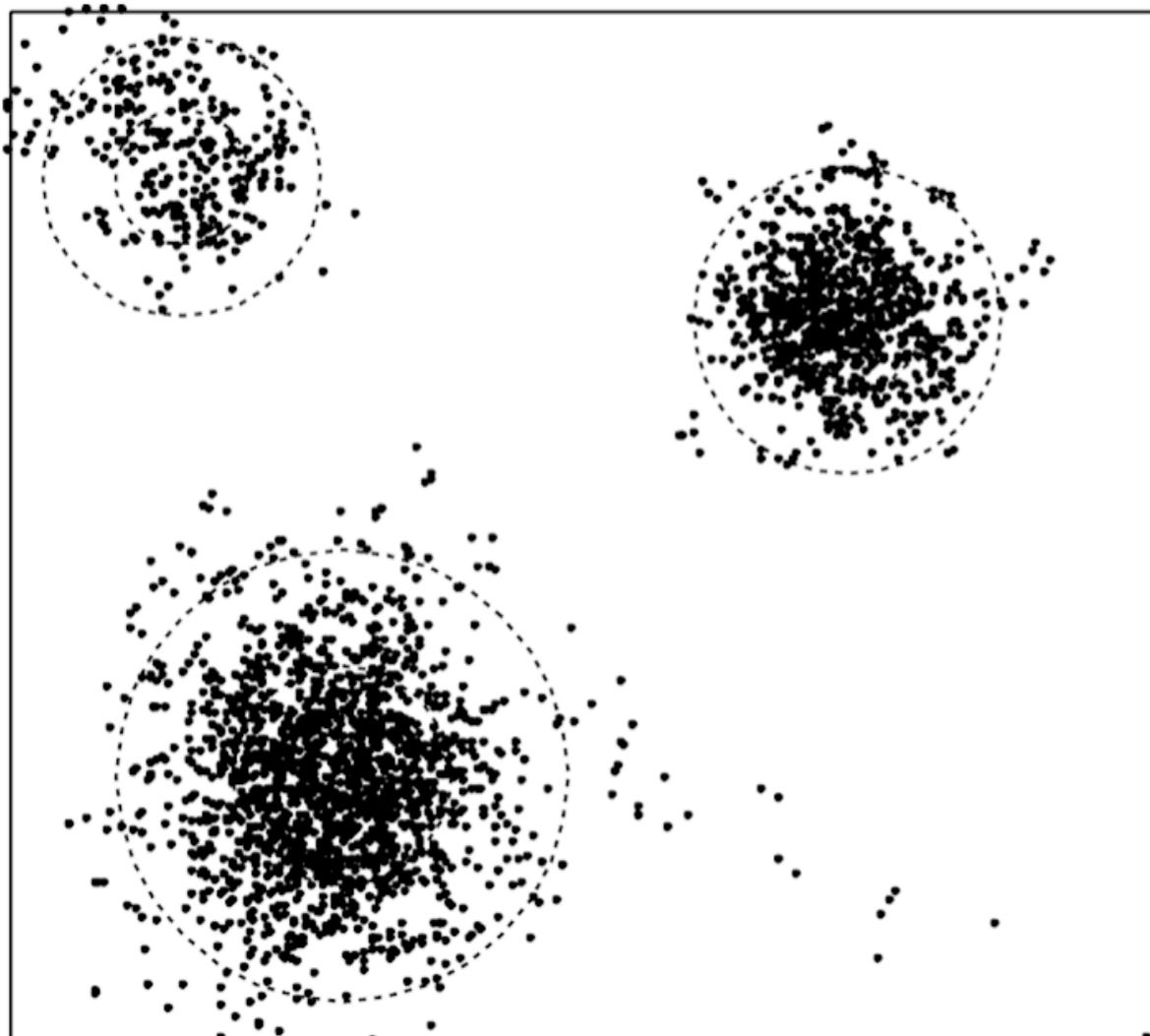


- Using MCMC rules the robot quickly finds one of the 3 hills
- First few steps are not representative of the distribution

First 100 steps by the robot

# MCMC approximation

(courtesy of Paul Lewis)



## How good is the approximation?

- 51.2% of points inside 50% contours
- 93.6% of points inside 95% contours

The more steps, the better the accuracy

5000 steps by the robot

# Operators

---

## Target distribution

- This is the **posterior** in BEAST2:  $P(\text{EvoSeq} | \text{ACAC...})$
- MCMC is a directed random walk through the state space and samples the target distribution
- How to pick the next state to sample?

## Proposal distribution

- Used to decide where to step to next
- The choice only affects the **efficiency** of the algorithm
- In BEAST1 and BEAST2 operators are used to propose the next step
- A parameter (or multiple parameters) are selected and perturbed to propose a step

# Operators

## Target

- This
- MC
- spa
- How
- Operators are a part of the MCMC **algorithm**,  
not the **model!**
- Tuning operators can help to improve efficiency,  
but should not change the results

## Proposal distribution

- Used to decide where to step to next
- The choice only affects the **efficiency** of the algorithm
- In BEAST1 and BEAST2 operators are used to propose the next step
- A parameter (or multiple parameters) are selected and perturbed to propose a step

# How many trees are there?

---

Genomes	Possible trees
4	15
5	105
6	945
7	10395
8	135135
9	2027025
10	34459425
20	$8.2 \times 10^{21}$
48	$3.21 \times 10^{70}$

$\approx$  The number of particles  
in the universe

# MCMC in practice

---

## Before

- Decide on the length of the chain  
(total number of steps to take)
- Decide on the sampling frequency  
(how often to record samples so  
that they are uncorrelated)

## After

- Discard burn-in  
(until stationary state is reached)
- Assess convergence and mixing
- Only then should we look at the estimates!

# MCMC in practice

---

## Before

- Decide on the length of the chain  
(total number of steps to take)
- Decide on the sampling frequency  
(how often to record samples so that they are uncorrelated)

## After

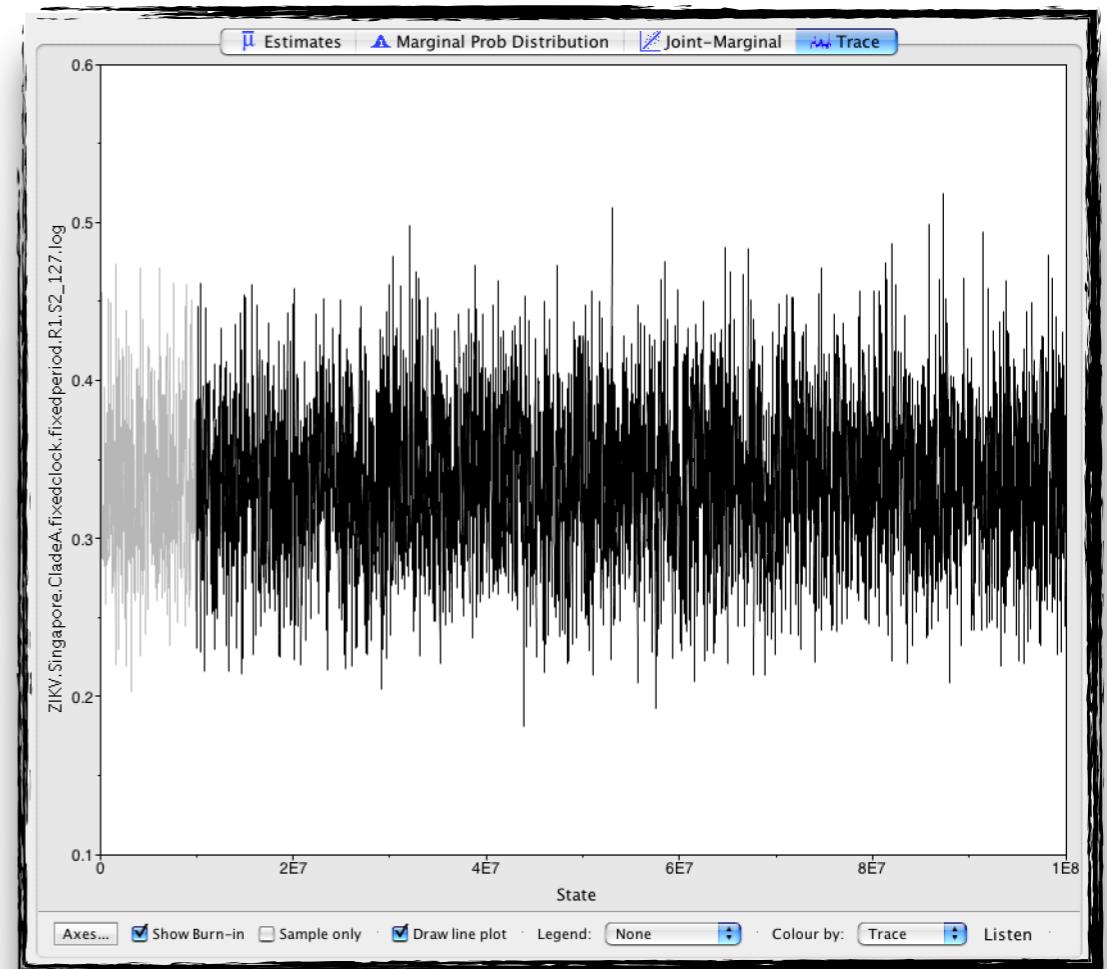
- Discard  
(up to the point where convergence is reached)
- Assess convergence and mixing
- Only then should we look at the estimates!

10,000 samples are plenty!  
(more is probably a waste of disk space...)

# What we hope will happen

---

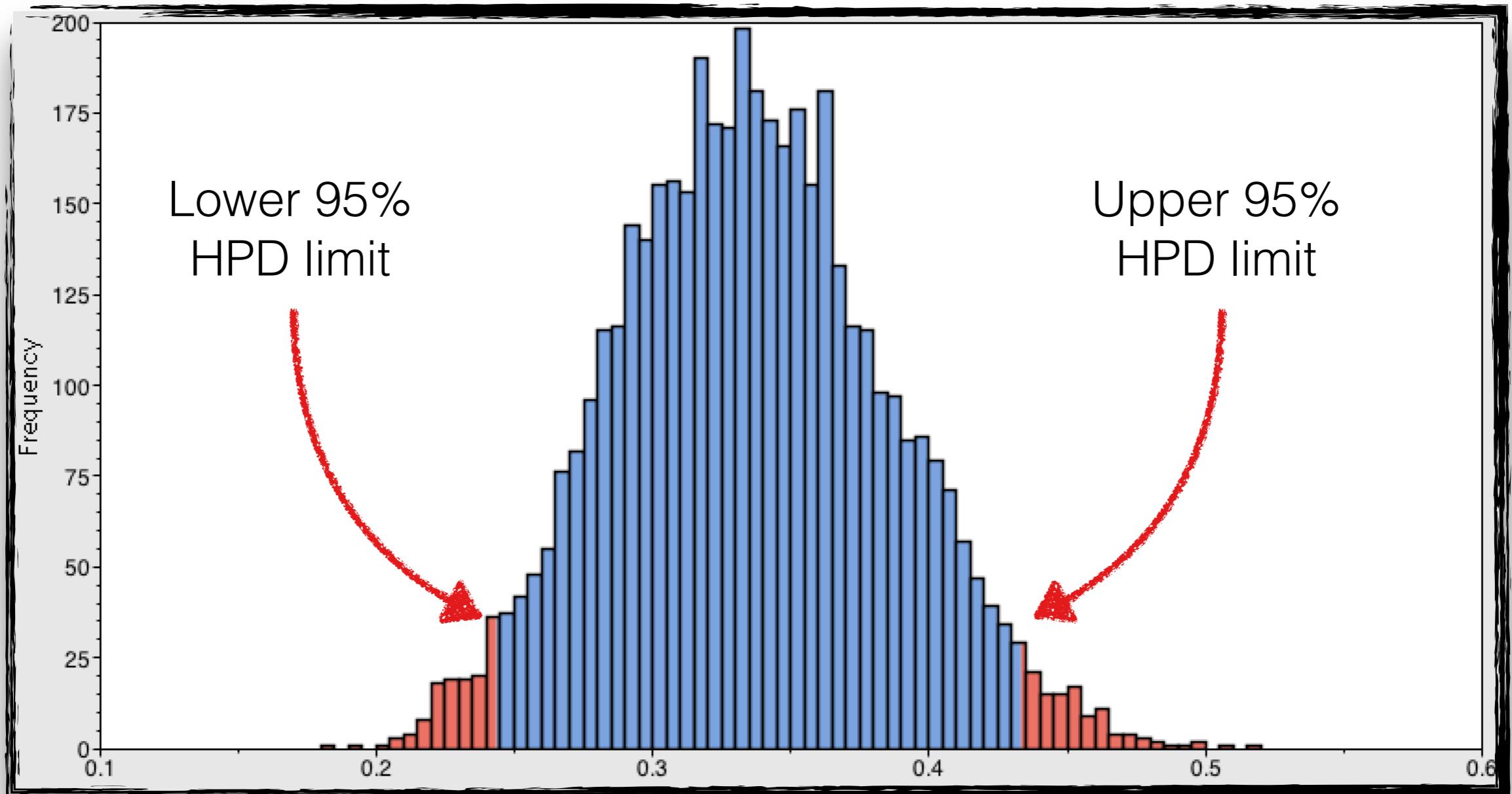
- The MCMC algorithm samples efficiently from high density areas of the posterior distribution
- We end up with a **good** approximation of the posterior distribution in **finite** time
- Appearance of white noise
- Everything is awesome!



Mixing well! 😊

# HPD intervals

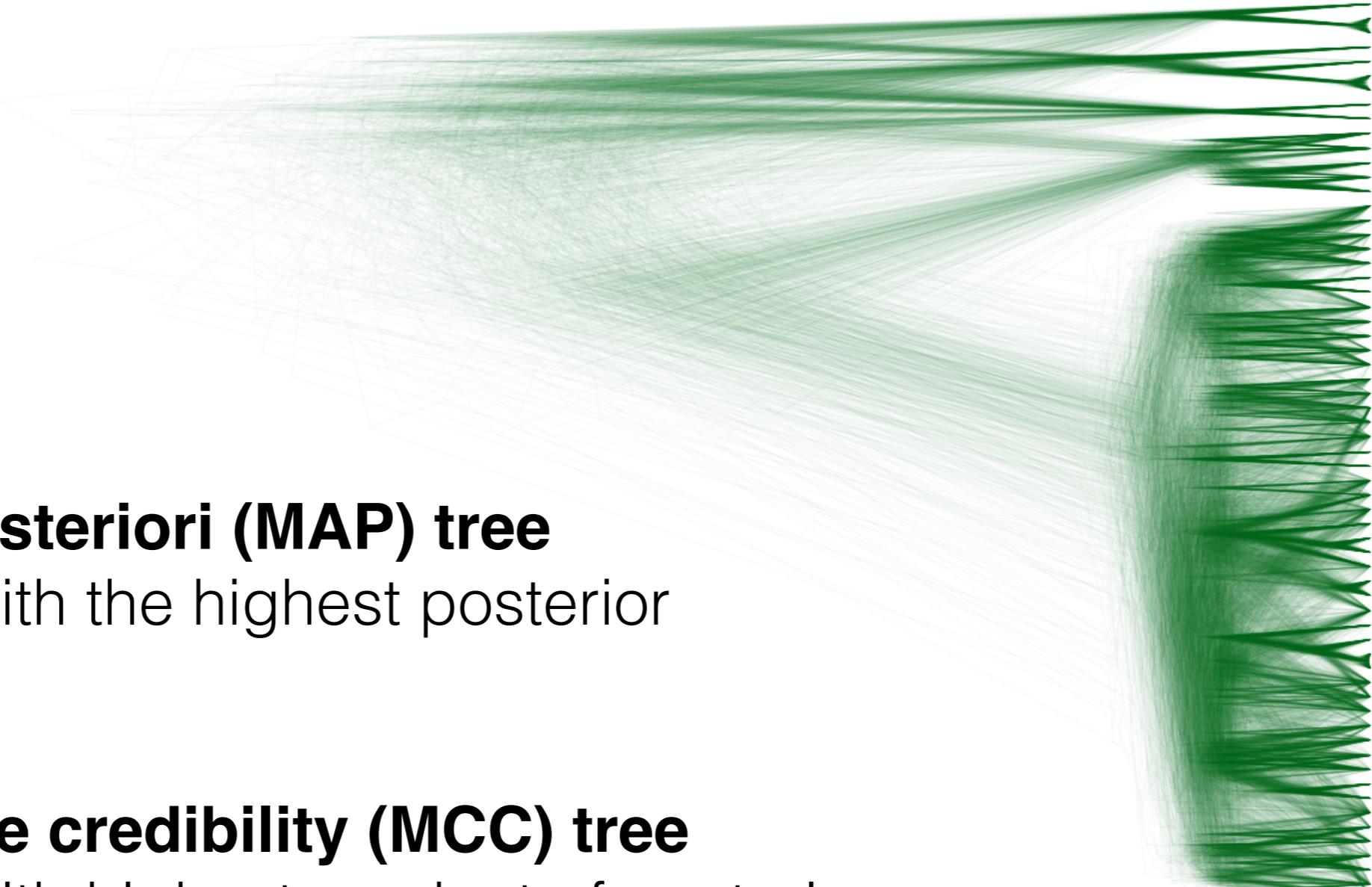
(Highest Posterior Density)



Smallest region that contains 95% of the posterior probability

# What about trees?

---



## **Maximum a posteriori (MAP) tree**

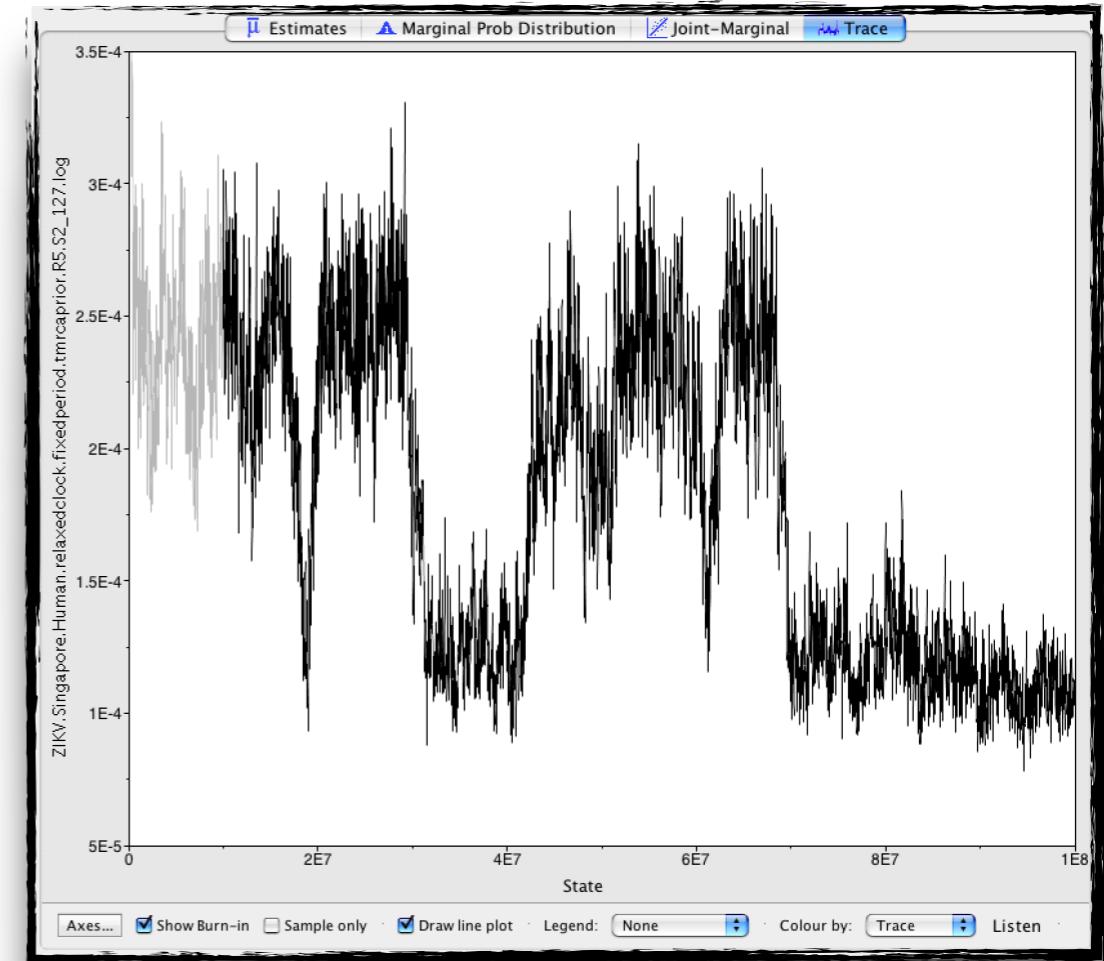
Sampled tree with the highest posterior probability

## **Maximum clade credibility (MCC) tree**

Sampled tree with highest product of posterior node probabilities

# Questions to ask...

- Is the chain **mixing** well?
- Are samples uniformly drawn from all over the stationary distribution?
- Do we have a “sticky chain?”



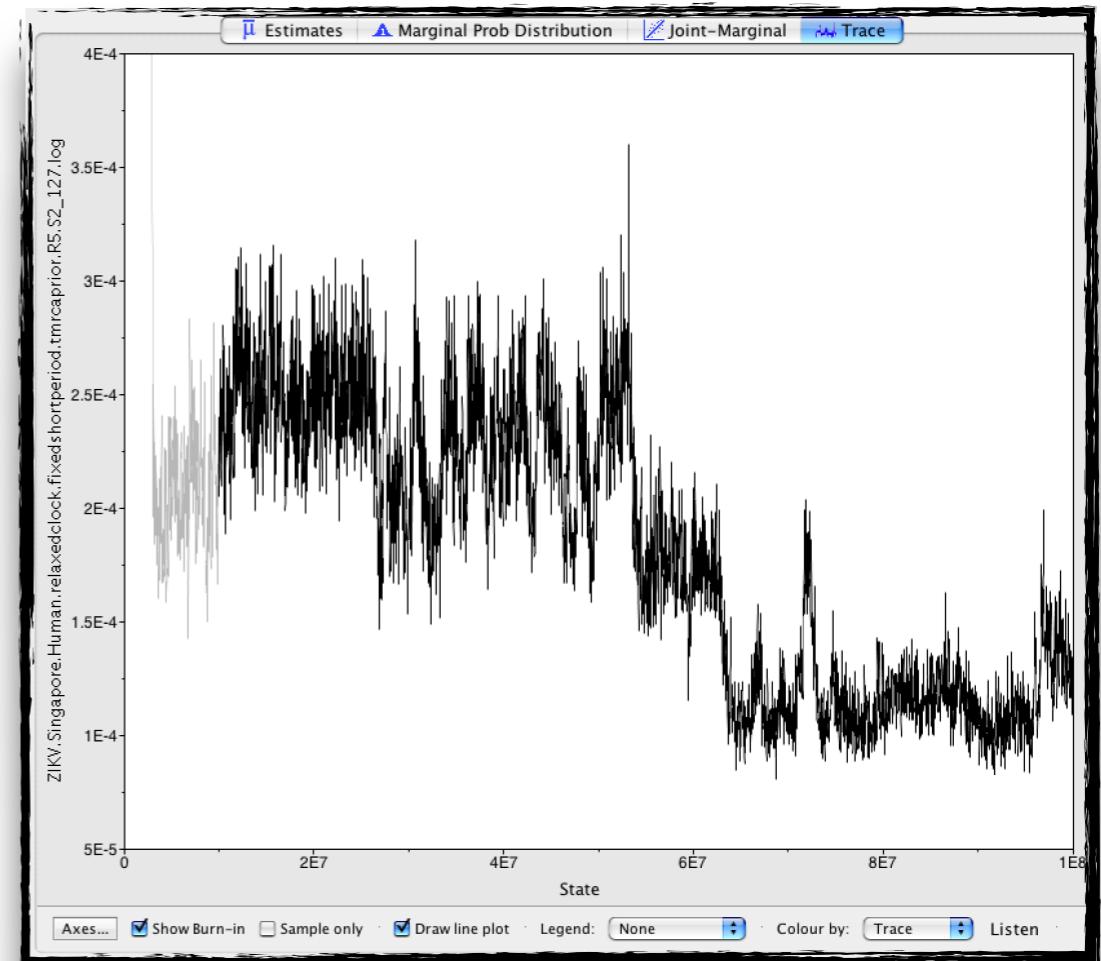
## Solutions

- MCMC gets stuck in some states for long times
- Tune operators to make better proposals

**Not mixing!** 😞

# Questions to ask...

- Has the chain **converged** to the stationary distribution?
- Did we pass the burn-in?



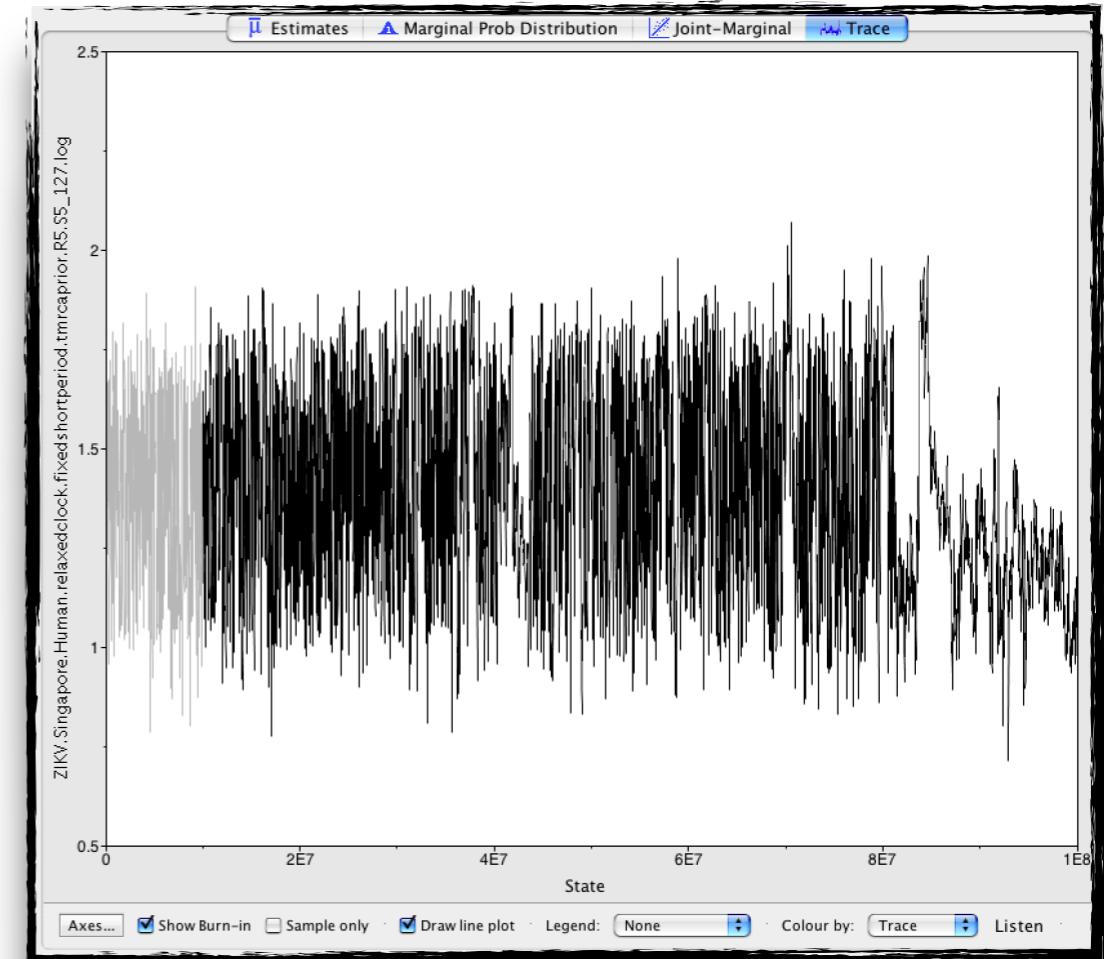
**Not converged!** 😓

**Solution:** Run for longer

# Questions to ask...

---

- Are we there yet?
- How do we know if the chain is long enough?



## Solution

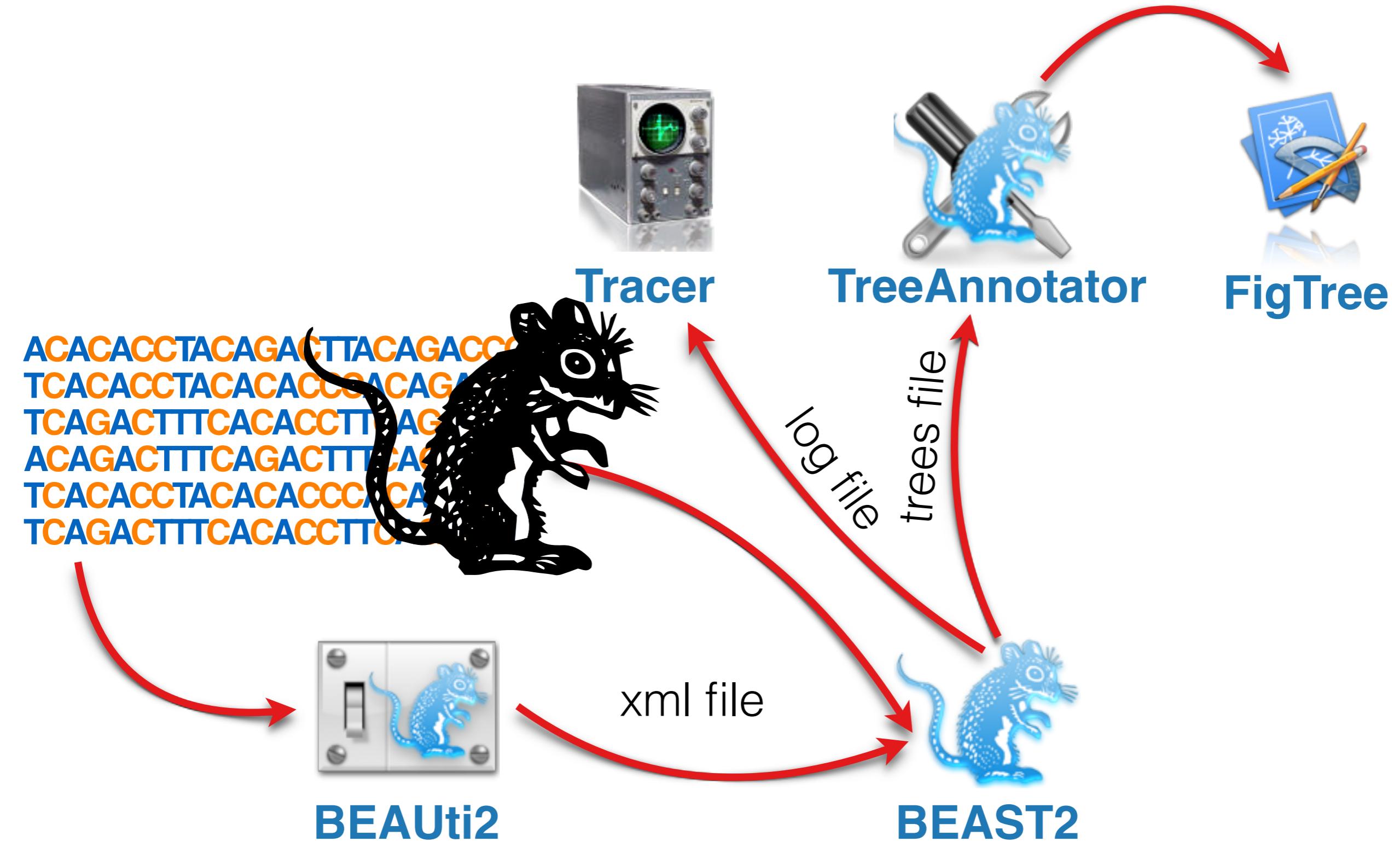
- Run multiple chains
- Combine chains
- Check that all chains give the same result

**Still not converged!** 😞

# FANTASTIC BEASTS

AND WHERE  
TO FIND THEM

# BEAST2 workflow



# BEAUti2

(<http://beast2.org>)



Graphical tool for setting up a BEAST2 analysis

## Input:

- Genetic sequence data
- Optional:
  - Sampling times
  - Sampling locations
  - Traits
  - etc.

## Output:

- Compact XML description of data, model and prior distributions that can be run in BEAST2

BEAUti 2: Standard /Users/louis/Documents/Taming\_the\_BEAST/Tutorials-Git/Introduction-to-BEAST2/xml/Primates.xml

Partitions Tip Dates Site Model Clock Model Priors MCMC

Link Site Models Unlink Site Models Link Clock Models Unlink Clock Models Link Trees Unlink Trees

Name	File	Taxa	Sites	Data Type	Site Model	Clock Model	Tree	...
noncoding	primate-mtDNA	12	205	nucleotide	noncoding	clock	tree	<input type="checkbox"/>
1stpos	primate-mtDNA	12	231	nucleotide	1stpos	clock	tree	<input type="checkbox"/>
2ndpos	primate-mtDNA	12	231	nucleotide	2ndpos	clock	tree	<input type="checkbox"/>
3rdpos	primate-mtDNA	12	231	nucleotide	3rdpos	clock	tree	<input type="checkbox"/>

+ - r Split

BEAUTi 2: Standard /Users/louis/Documents/Taming\_the\_BEAST/Tutorials-Git/Introduction-to-BEAST2/xml/Primates.xml

Partitions Tip Dates Site Model Clock Model Priors MCMC

▶ Tree.t:tree      Calibrated Yule Model

▶ birthRateY.t:tree      Gamma      initial = [1.0]  $[-\infty, \infty]$       Calibrated Yule speciation process birth rate for t:3rdpos

▶ clockRate.c:clock      Uniform      initial = [1.0]  $[-\infty, \infty]$       substitution rate of partition c:3rdpos

▶ gammaShape.s:1stpos      Exponential      initial = [1.0]  $[-\infty, \infty]$       Prior on gamma shape for partition s:1stpos

▶ gammaShape.s:2ndpos      Exponential      initial = [1.0]  $[-\infty, \infty]$       Prior on gamma shape for partition s:2ndpos

▶ gammaShape.s:3rdpos      Exponential      initial = [1.0]  $[-\infty, \infty]$       Prior on gamma shape for partition s:3rdpos

▶ gammaShape.s:noncoding      Exponential      initial = [1.0]  $[-\infty, \infty]$       Prior on gamma shape for partition s:noncoding

▶ kappa.s:1stpos      Log Normal      initial = [2.0]  $[0.0, \infty]$       HKY transition-transversion parameter of partition s:1stpos

▶ kappa.s:2ndpos      Log Normal      initial = [2.0]  $[0.0, \infty]$       HKY transition-transversion parameter of partition s:2ndpos

▶ kappa.s:3rdpos      Log Normal      initial = [2.0]  $[0.0, \infty]$       HKY transition-transversion parameter of partition s:3rdpos

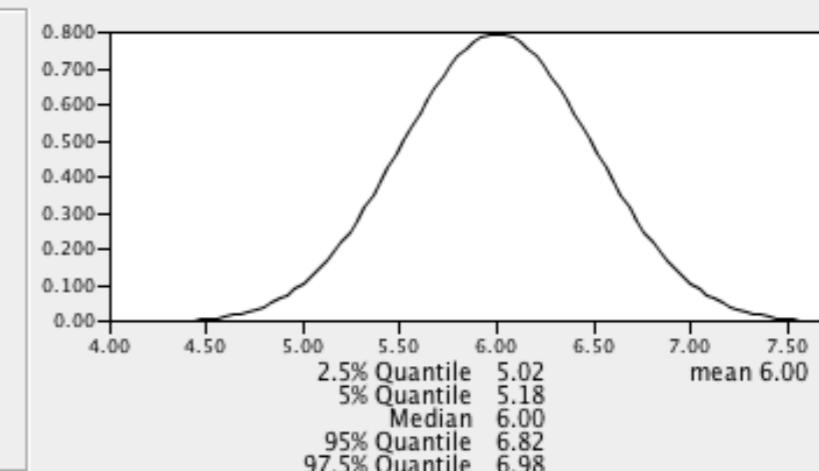
▶ kappa.s:noncoding      Log Normal      initial = [2.0]  $[0.0, \infty]$       HKY transition-transversion parameter of partition s:noncoding

▶ human-chimp.prior      Normal      initial = [6.0]  $[4.5, 7.5]$       monophyletic

Mean: 6.0       estimate

Sigma: 0.5       estimate

Offset: 0.0



2.5% Quantile 5.02  
 5% Quantile 5.18  
 Median 6.00  
 95% Quantile 6.82  
 97.5% Quantile 6.98

Tipsonly  
 Use Originate

Primates\_long.xml UNREGISTERED

```
39 <run id="mcmc" spec="MCMC" chainLength="2500000">
40   <state id="state" storeEvery="5000">
41     <tree id="Tree.t:tree" name="stateNode">
42       <taxonset id="TaxonSet.noncoding" spec="TaxonSet">
43         <alignment id="noncoding" spec="FilteredAlignment" filter="1,458-659,897-898">
44           <data idref="primate-mtDNA"/>
45         </alignment>
46       </taxonset>
47     </tree>
48     <parameter id="mutationRate.s:noncoding" name="stateNode">1.0</parameter>
49     <parameter id="gammaShape.s:noncoding" name="stateNode">1.0</parameter>
50     <parameter id="kappa.s:noncoding" lower="0.0" name="stateNode">2.0</parameter>
51     <parameter id="kappa.s:1stpos" lower="0.0" name="stateNode">2.0</parameter>
52     <parameter id="gammaShape.s:1stpos" name="stateNode">1.0</parameter>
53     <parameter id="mutationRate.s:1stpos" name="stateNode">1.0</parameter>
54     <parameter id="kappa.s:2ndpos" lower="0.0" name="stateNode">2.0</parameter>
55     <parameter id="gammaShape.s:2ndpos" name="stateNode">1.0</parameter>
56     <parameter id="mutationRate.s:2ndpos" name="stateNode">1.0</parameter>
57     <parameter id="kappa.s:3rdpos" lower="0.0" name="stateNode">2.0</parameter>
58     <parameter id="gammaShape.s:3rdpos" name="stateNode">1.0</parameter>
59     <parameter id="mutationRate.s:3rdpos" name="stateNode">1.0</parameter>
60     <parameter id="birthRateY.t:tree" name="stateNode">1.0</parameter>
61     <parameter id="clockRate.c:clock" name="stateNode">1.0</parameter>
62   </state>
63
64
65   <init id="RandomTree.t:tree" spec="beast.evolution.tree.RandomTree" estimate="false" initial="@Tree.t:tree" taxa="@noncoding">
66     <populationModel id="ConstantPopulation0.t:tree" spec="ConstantPopulation">
67       <parameter id="randomPopSize.t:tree" name="popSize">1.0</parameter>
68     </populationModel>
69   </init>
70
71   <distribution id="posterior" spec="util.CompoundDistribution">
72     <distribution id="prior" spec="util.CompoundDistribution">
73       <distribution id="CalibratedYuleModel.t:tree" spec="beast.evolution.speciation.CalibratedYuleModel" birthRate="@birthRateY.t:tree" tree="@Tree.t:tree"/>
74       <prior id="CalibratedYuleBirthRatePrior.t:tree" name="distribution" x="@birthRateY.t:tree">
75         <Gamma id="Gamma.0" name="distr">
76           <parameter id="RealParameter.0" estimate="false" name="alpha">0.001</parameter>
77           <parameter id="RealParameter.01" estimate="false" name="beta">1000.0</parameter>
78         </Gamma>
79       </prior>
80       <prior id="ClockPrior.c:clock" name="distribution" x="@clockRate.c:clock">
81         <Uniform id="Uniform.0" name="distr" upper="Infinity"/>
82       </prior>
83     </distribution>
84   </distribution>
85 
```

# BEAST2

(<http://beast2.org>)



- Bayesian **e**volutionary **a**nalysis by **s**ampling **t**rees
- Performs MCMC analyses of sequences under selected sequence evolution and tree model
- Similar to BEAST 1.10 but completely separate and generally incompatible
- BEAST2 and BEAST1 have a common origin, have much of the same functionality but have diverged over time
- BEAST2 has a modular design that makes it easy to extend
- GUI interface but can also be run from command line (e.g. on a cluster)

## Input:

- xml model description file

## Outputs:

- log file
- trees file
- state file

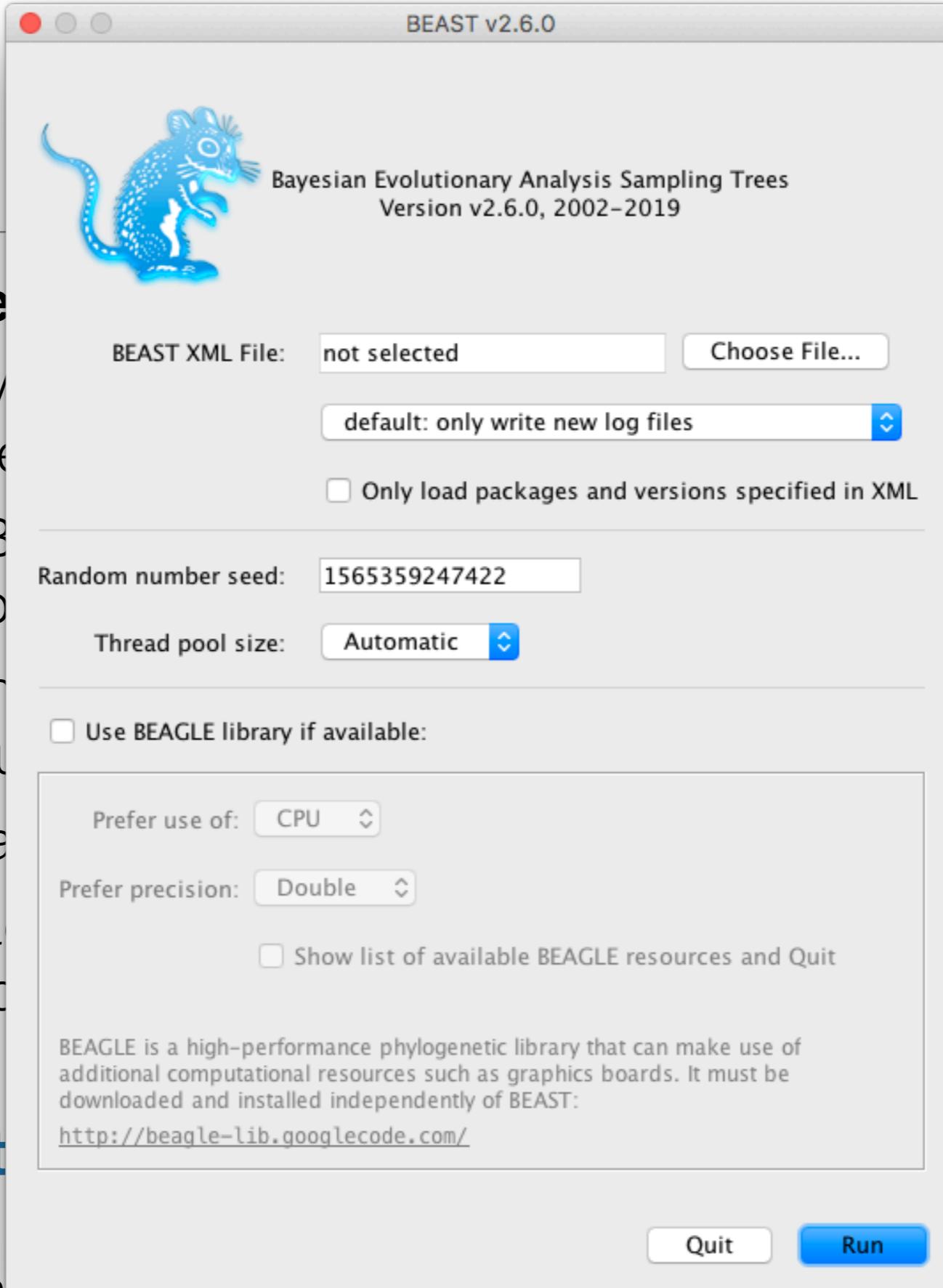
# BEAST2

(<http://beast2.org>)

- Bayesian estimation
- Performs MCMC analysis on sequence evolution
- Similar to BEAGLE, but incompatible
- BEAST2 and BEAGLE have the same functions
- BEAST2 has a GUI
- GUI interface (e.g. on a desktop)

## Input

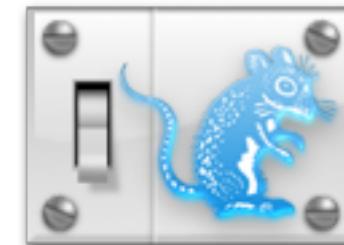
- XML description file



ected  
d generally  
e much of  
to extend  
ine

- state file

# BEAST2 packages



- BEAST2 is organised into a central "core" together with a large number of separate packages
- Packages can be developed by anybody — **including you!**
- Can be directly integrated into BEAST2 and updated frequently without waiting for a full BEAST2 release
- Packages add new models or completely new functionality
  - phylogeography
  - bacterial ARG inference
  - morphological models
  - multispecies coalescent
  - model selection and averaging
  - stochastic simulations
  - linguistic analyses
  - ...
- Install new packages through the package manager in BEAUti

BEAST 2 Package Manager

List of available packages for BEAST v2.6.\*

Name	Installed	Latest	Dependencies	Link	Detail
BEAST	2.6.0	2.6.1		<a href="#">[ ]</a>	<b>BEAST core</b>
Babel		0.2.1	BEASTLabs	<a href="#">[ ]</a>	BABEL = BEAST analysis backing effective linguistics
bacter	2.2.3	2.2.3		<a href="#">[ ]</a>	Bacterial ARG inference.
BADTRIP		1.0.0		<a href="#">[ ]</a>	Infer transmission time for non-haplotype data and epi data
BASTA		3.0.1		<a href="#">[ ]</a>	Bayesian structured coalescent approximation
bdmm	0.3.5	0.3.5	MultiTypeTree, MASTER	<a href="#">[ ]</a>	pre-release of multitype birth-death model (aka birth-death-migration model)
BDSKY	1.4.5	1.4.5		<a href="#">[ ]</a>	birth death skyline – handles serially sampled tips, piecewise constant rate changes through time and sampled ancestors
BEAST_CLASSIC		1.5.0	BEASTLabs	<a href="#">[ ]</a>	BEAST classes ported from BEAST 1 in wrappers
BEASTLabs	1.9.0	1.9.0		<a href="#">[ ]</a>	BEAST utilities, such as Script, multi monophyletic constraints
BEASTvntr		0.1.3		<a href="#">[ ]</a>	Variable Number of Tandem Repeat data, such as microsatellites
Beasy		0.0.2	BEASTLabs	<a href="#">[ ]</a>	Makes it easier to construct models: Automatic methods text generator, Beasy XML generator, and more
bModelTest	1.2.1	1.2.1	BEASTLabs	<a href="#">[ ]</a>	Bayesian model test for nucleotide subst models, gamma rate heterogeneity and invariant sites
BREAK_AWAY		1.0.0	BEASTLabs, GEO_SPHERE	<a href="#">[ ]</a>	break-away model of phylogeography
CA		2.0.0		<a href="#">[ ]</a>	Bayesian estimation of clade ages based on probabilities of fossil sampling
CoalRe		0.0.4		<a href="#">[ ]</a>	Infer viral reassortment networks
CodonSubstModels		1.1.3		<a href="#">[ ]</a>	Codon substitution models
CoupledMCMC		0.1.7	BEASTLabs	<a href="#">[ ]</a>	Coupled MCMC (parallel Tempering or MC3)
DENIM		1.0.0		<a href="#">[ ]</a>	Divergence Estimation Notwithstanding ILS and Migration
EpiInf		7.1.5	SA	<a href="#">[ ]</a>	BD/SIR/SIS epidemic trajectory inference.
FLC		1.1.0		<a href="#">[ ]</a>	Flexible local clock model
GEO_SPHERE		1.3.0	BEASTLabs	<a href="#">[ ]</a>	Whole world phylogeography
Mascot		1.2.2		<a href="#">[ ]</a>	Marginal approximation of the structured coalescent
MASTER	6.1.1	6.1.1		<a href="#">[ ]</a>	Stochastic population dynamics simulation
MGSM		0.3.0		<a href="#">[ ]</a>	Multi-gamma and relaxed gamma site models
MM		1.1.1		<a href="#">[ ]</a>	Enables models of morphological character evolution
MODEL_SELECTION		1.5.1	BEASTLabs	<a href="#">[ ]</a>	Select models through path sampling/stepping stone analysis
MSBD		1.1.0		<a href="#">[ ]</a>	Multi-state birth-death prior with state-specific birth and death rates
MultiTypeTree	7.0.1	7.0.1		<a href="#">[ ]</a>	Structured coalescent inference
NS		1.1.0	MODEL_SELECTION, BEASTLabs	<a href="#">[ ]</a>	Nested sampling for model selection and posterior inference
PhyDyn		1.3.4		<a href="#">[ ]</a>	PhyDyn: Epidemiological modelling with BEAST
phylodynamics		1.3.0	BDSKY	<a href="#">[ ]</a>	BDSIR and Stochastic Coalescent
PoMo		1.0.1		<a href="#">[ ]</a>	PoMo, a substitution model that separates mutation and drift processes
SA	2.0.2	2.0.2	BEASTLabs	<a href="#">[ ]</a>	Sampled ancestor trees
SCOTTI		2.0.1		<a href="#">[ ]</a>	Structured COalescent Transmission Tree Inference
SNAPP		1.5.0		<a href="#">[ ]</a>	SNP and AFLP Phylogenies
SpeciesNetwork		0.12.2		<a href="#">[ ]</a>	Multispecies network coalescent (MSNC) inference of introgression and hybridization
SSM		1.1.0		<a href="#">[ ]</a>	Standard Nucleotide Substitution Models
STACEY		1.2.5		<a href="#">[ ]</a>	Species delimitation and species tree estimation
StarBEAST2		0.15.5	SA, MM	<a href="#">[ ]</a>	Multispecies coalescent inference using multi-locus and fossil data
substBMA		1.2.3		<a href="#">[ ]</a>	Substitution Bayesian Model Averaging
TMA		1.0.0	MASTER, BEASTLabs, phylodynamics, BDSKY, TreeStat2	<a href="#">[ ]</a>	Tree model adequacy: test whether the tree prior used is adequate for your data
TreeStat2		0.0.2		<a href="#">[ ]</a>	Utility for calculating tree statistics from tree log file

 Latest

Install/Upgrade

Uninstall

Package repositories

Close

?

- Install new packages through the package manager in BEAUTi

# Tracer

(<http://beast.community>)



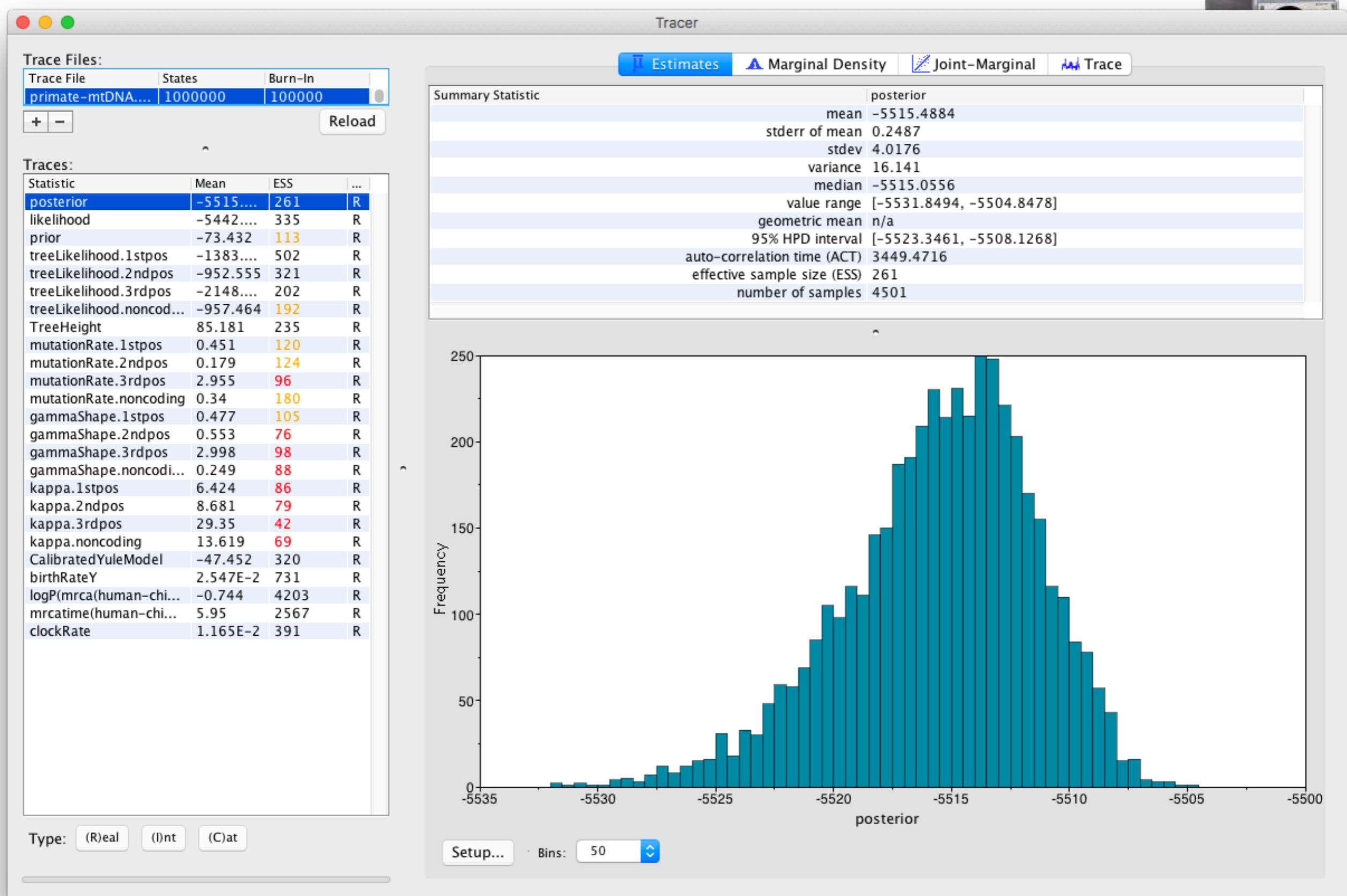
- Analyse (parameter) log files from BEAST2 runs
- Check mixing, ESS, ACT, parameter correlations
- Provides overview of posterior parameter estimates
- Easily compare several analyses
- Demographic reconstruction for some models (e.g. Bayesian Skyline Plot)
- Tracer is **primarily** a diagnostic tool — usually perform final analyses in a statistical package like R!

## Input:

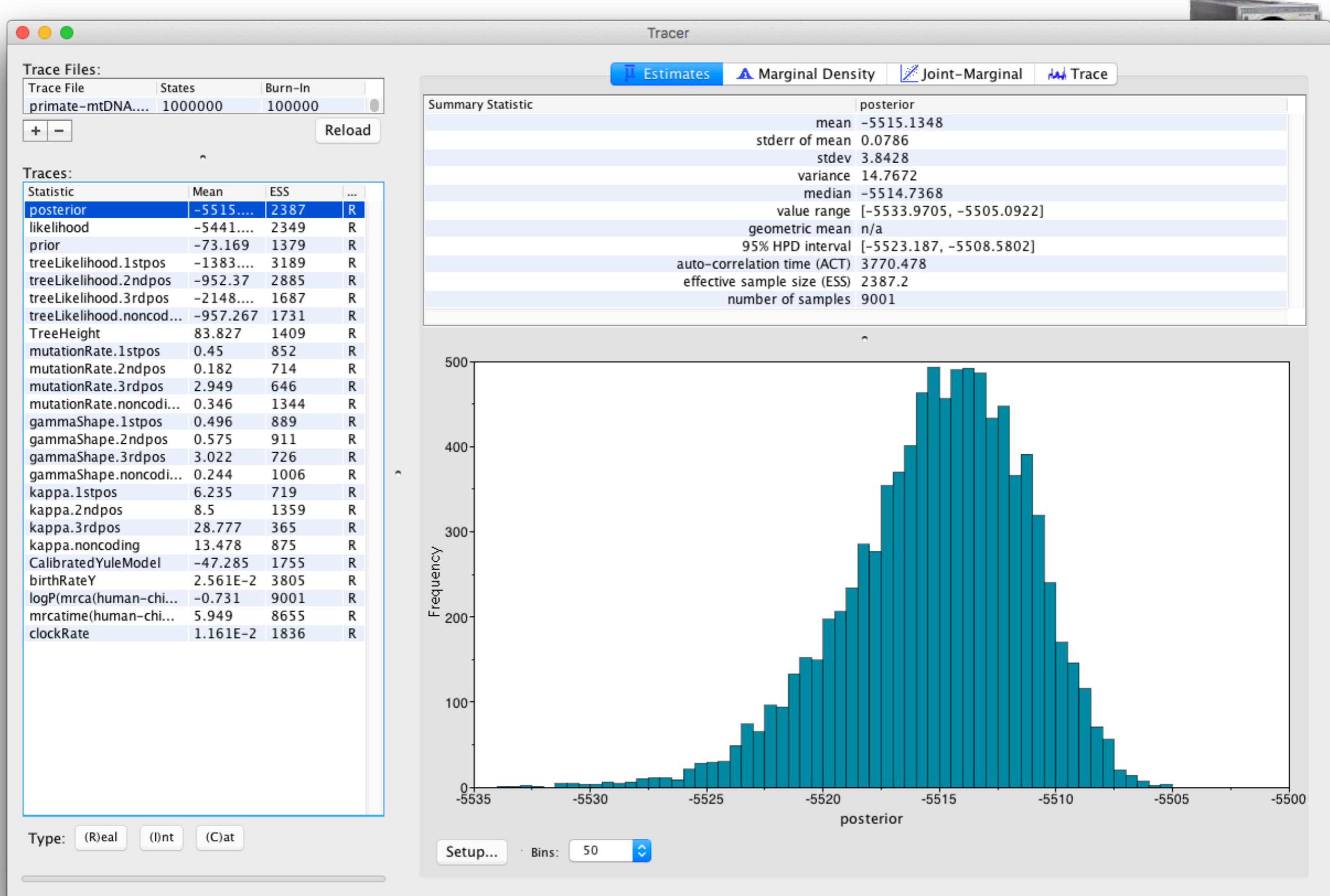
- log file

## Output:

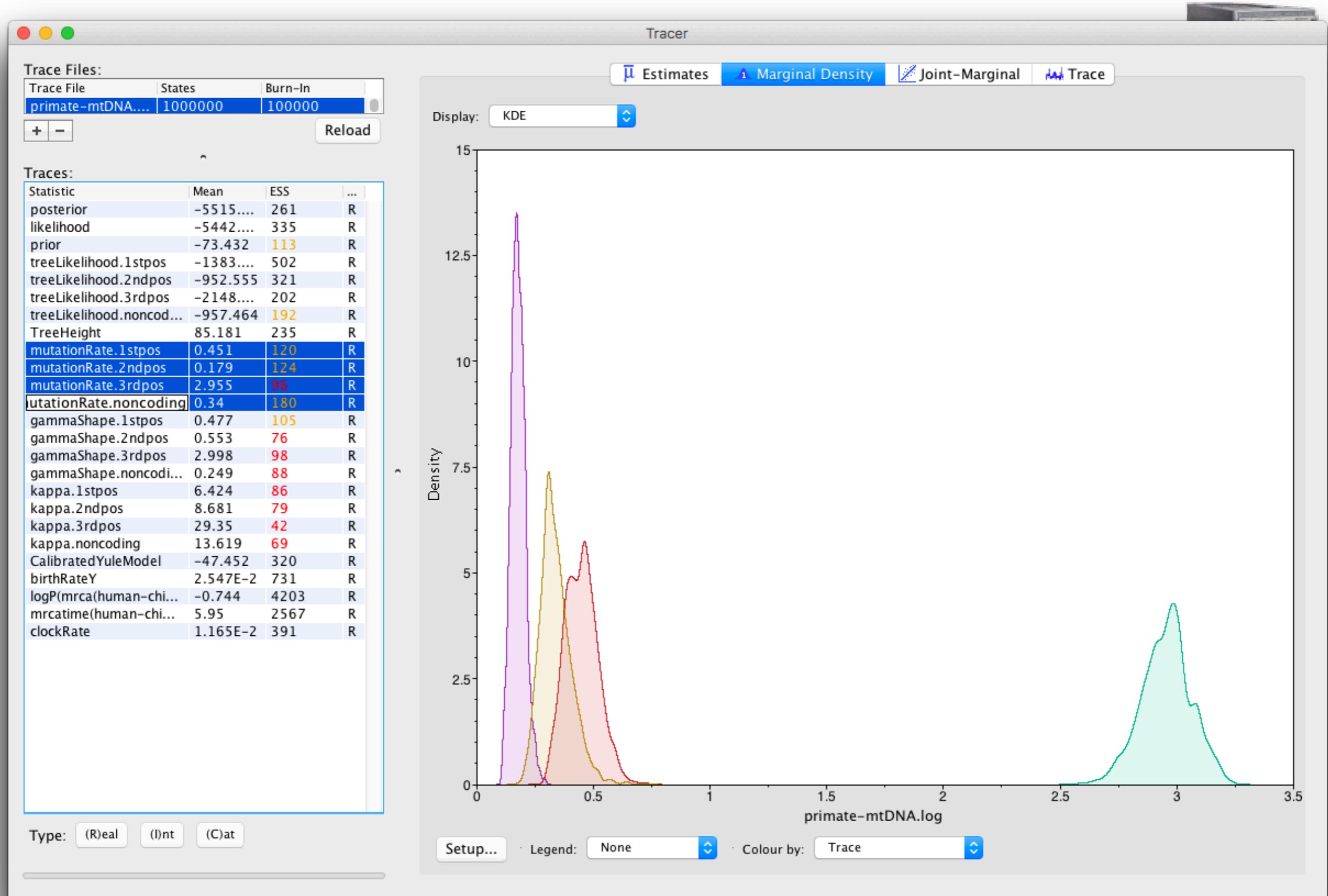
- Gain insight
- Demographic reconstructions



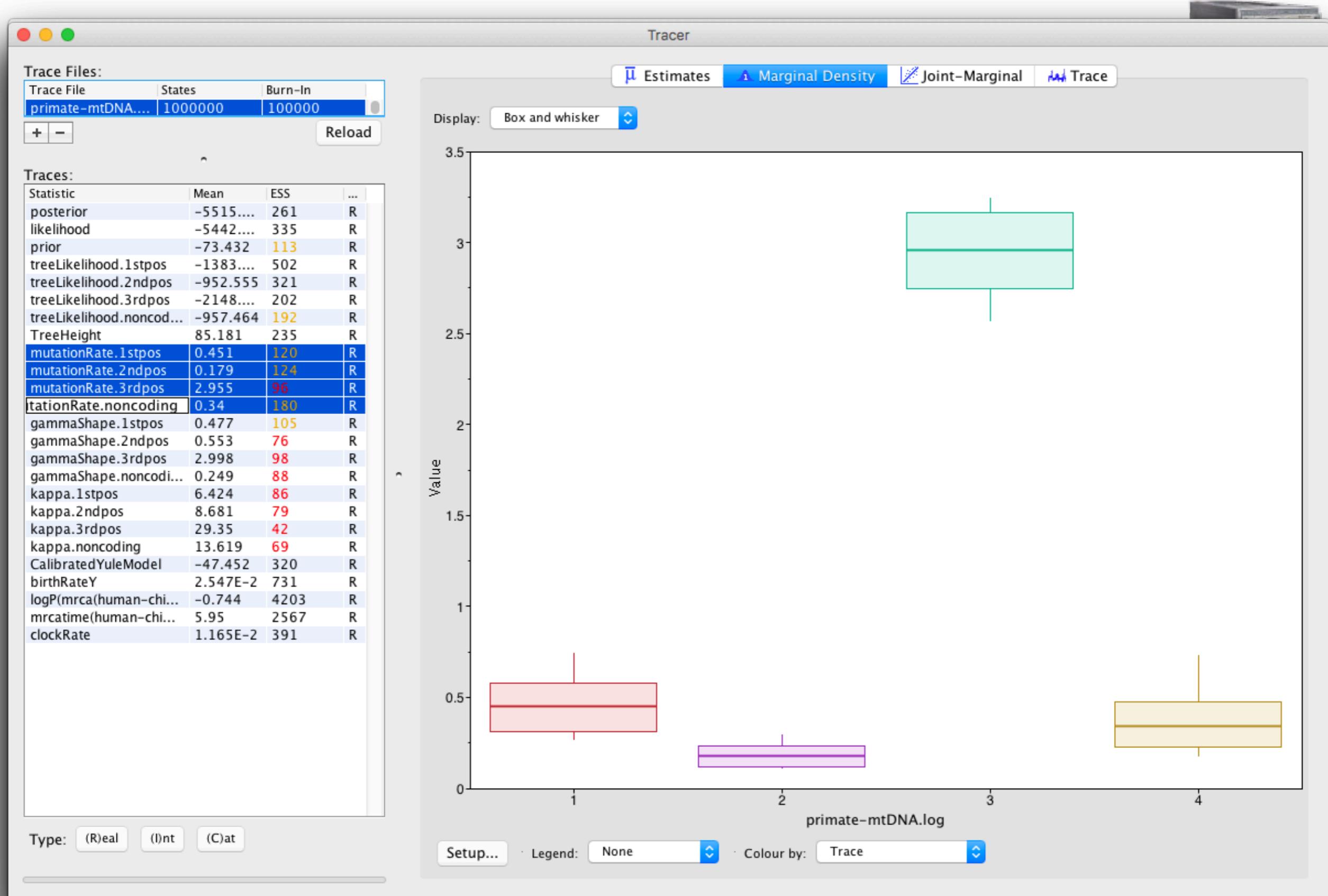
- Demographic reconstructions



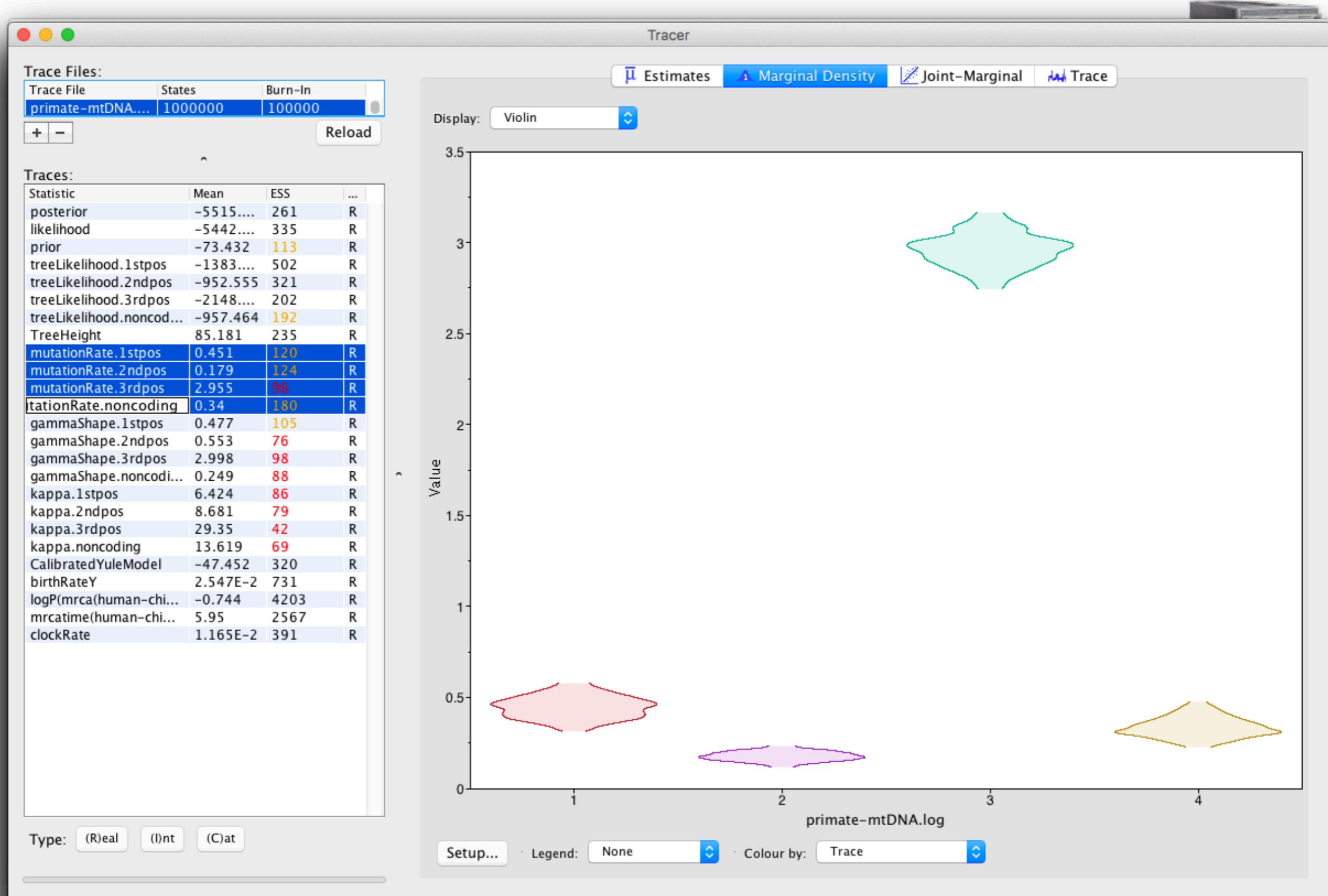
- Demographic reconstructions



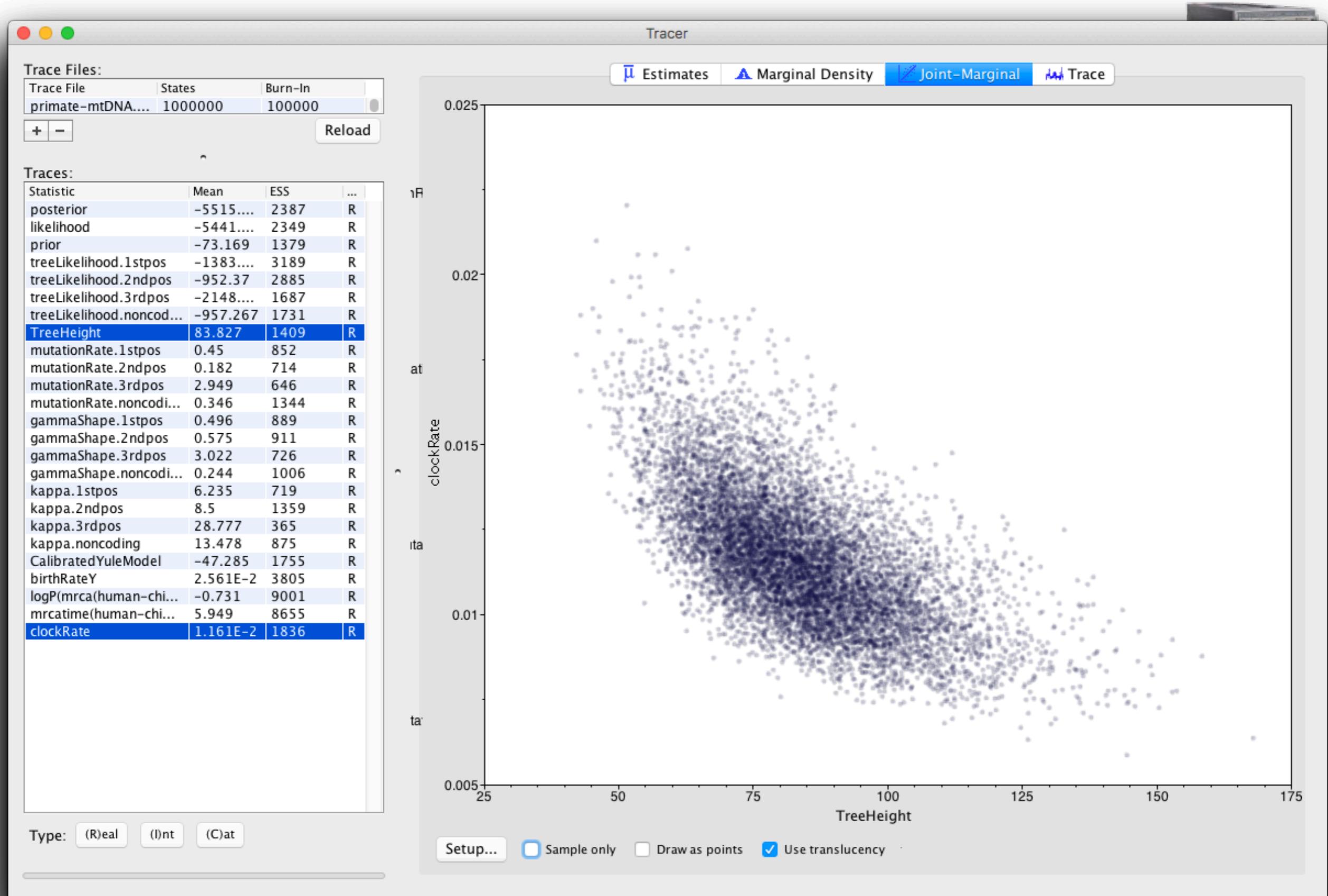
- Demographic reconstructions



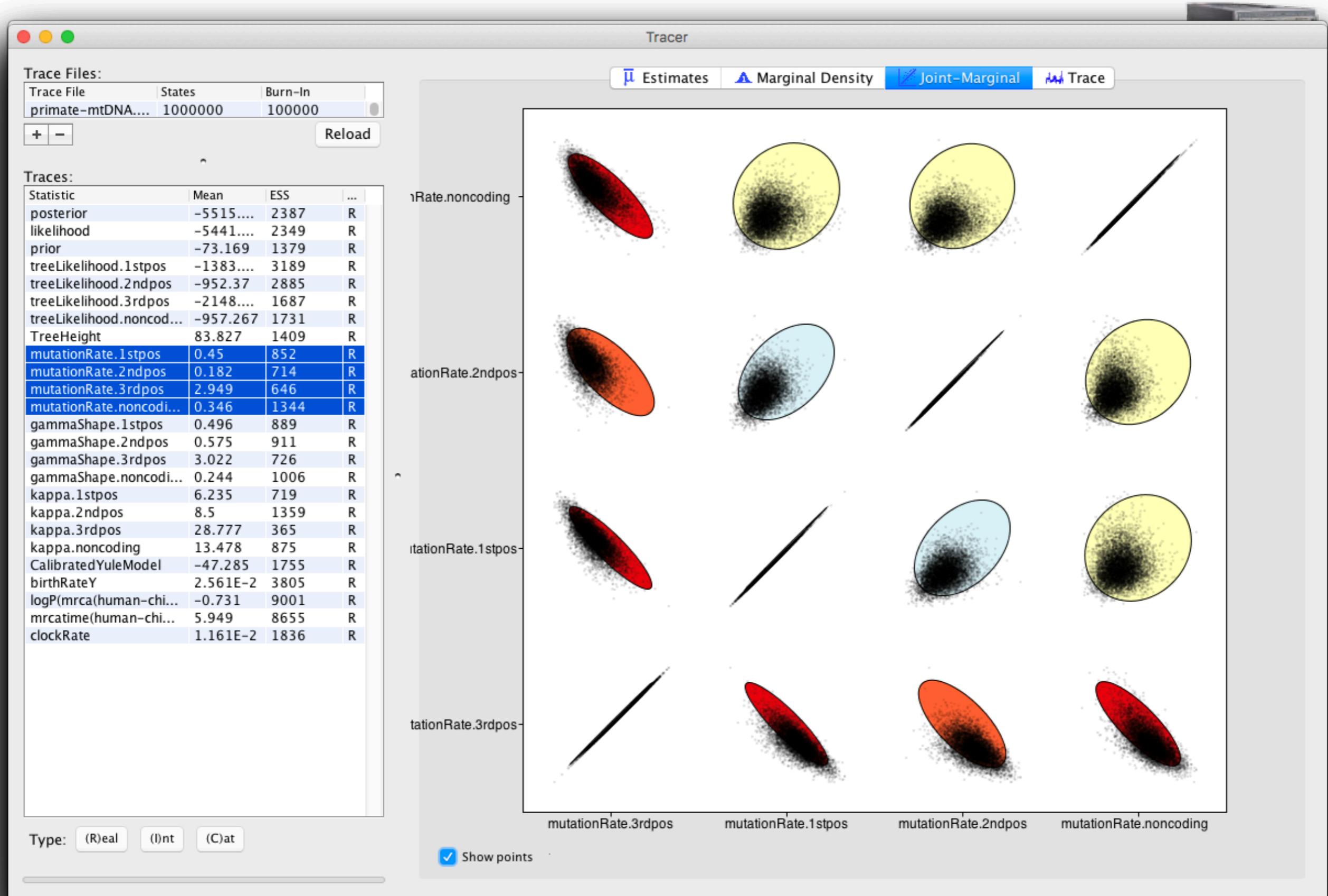
- Demographic reconstructions



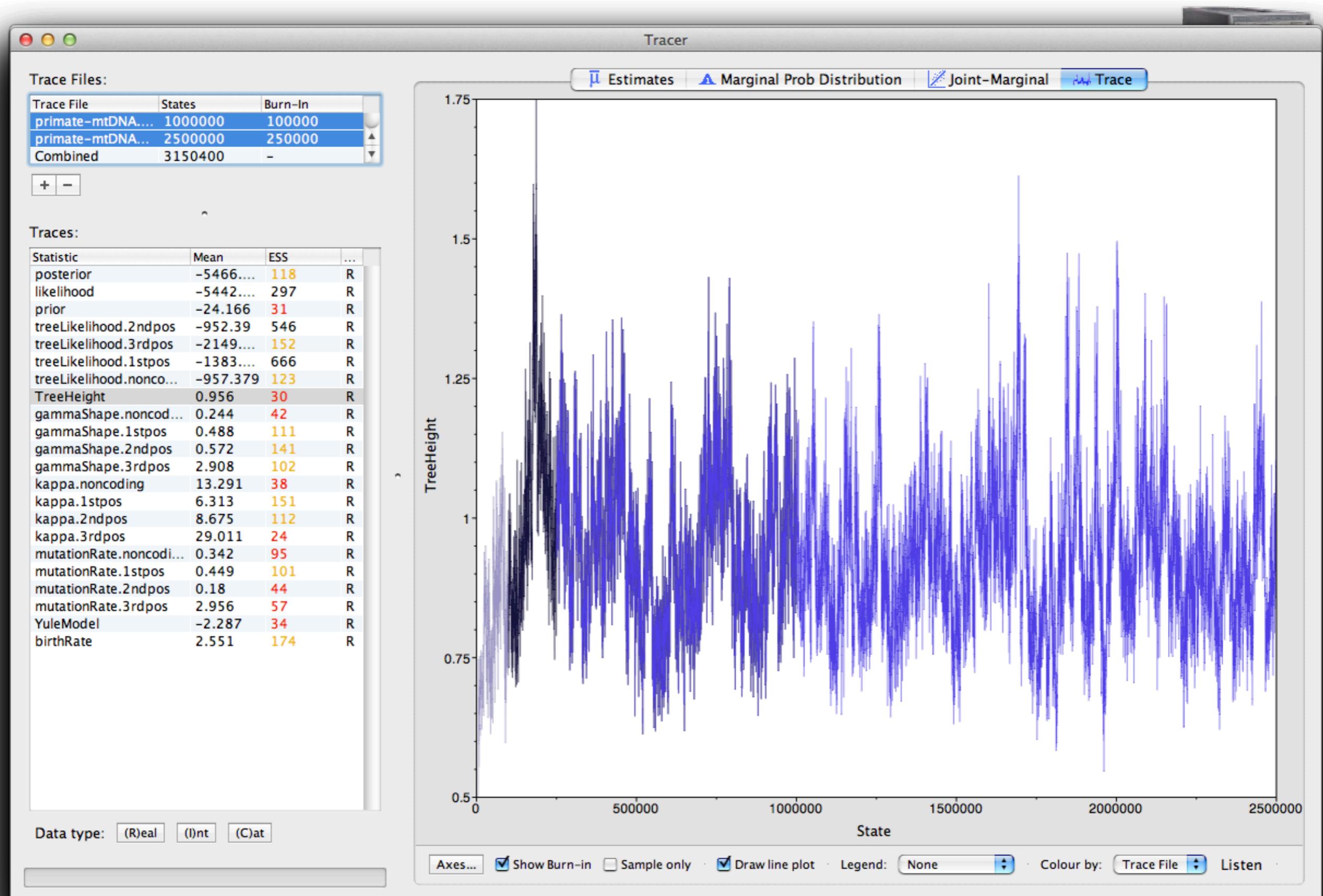
- Demographic reconstructions



- Demographic reconstructions

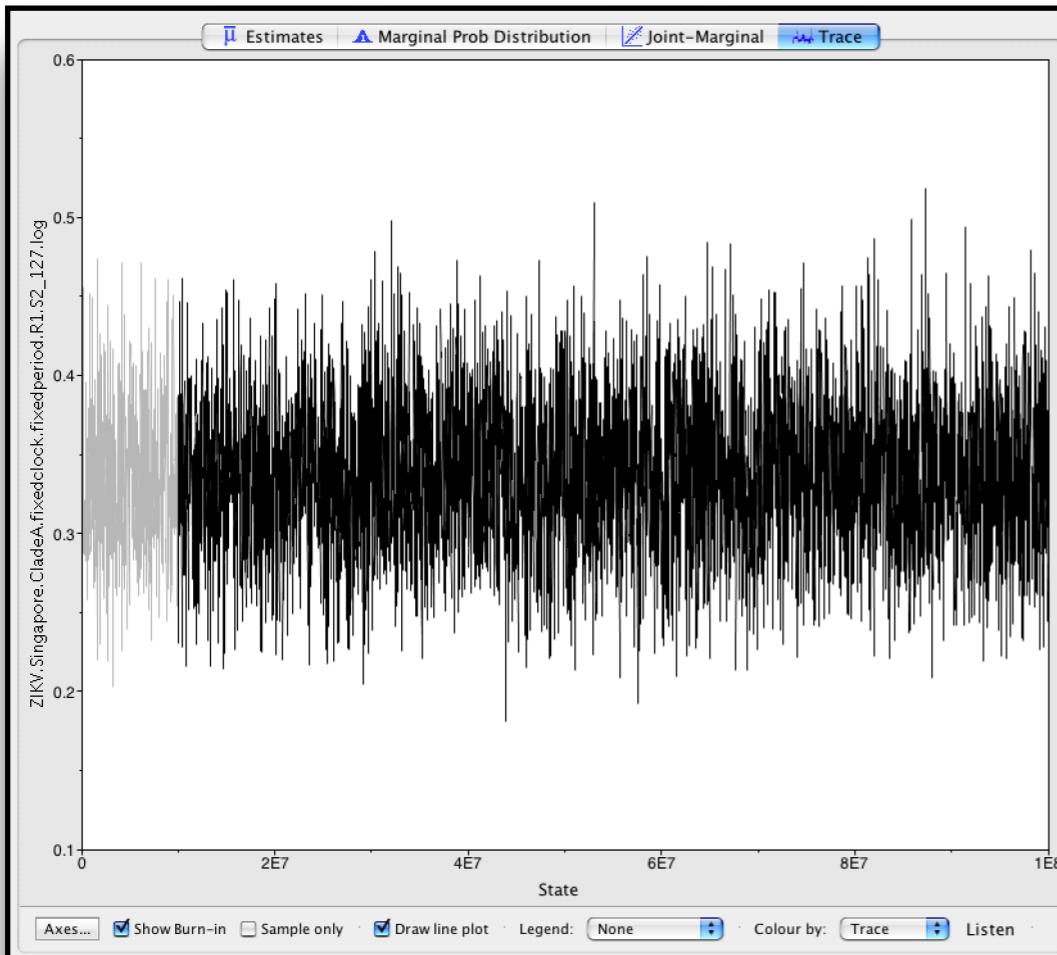


- Demographic reconstructions

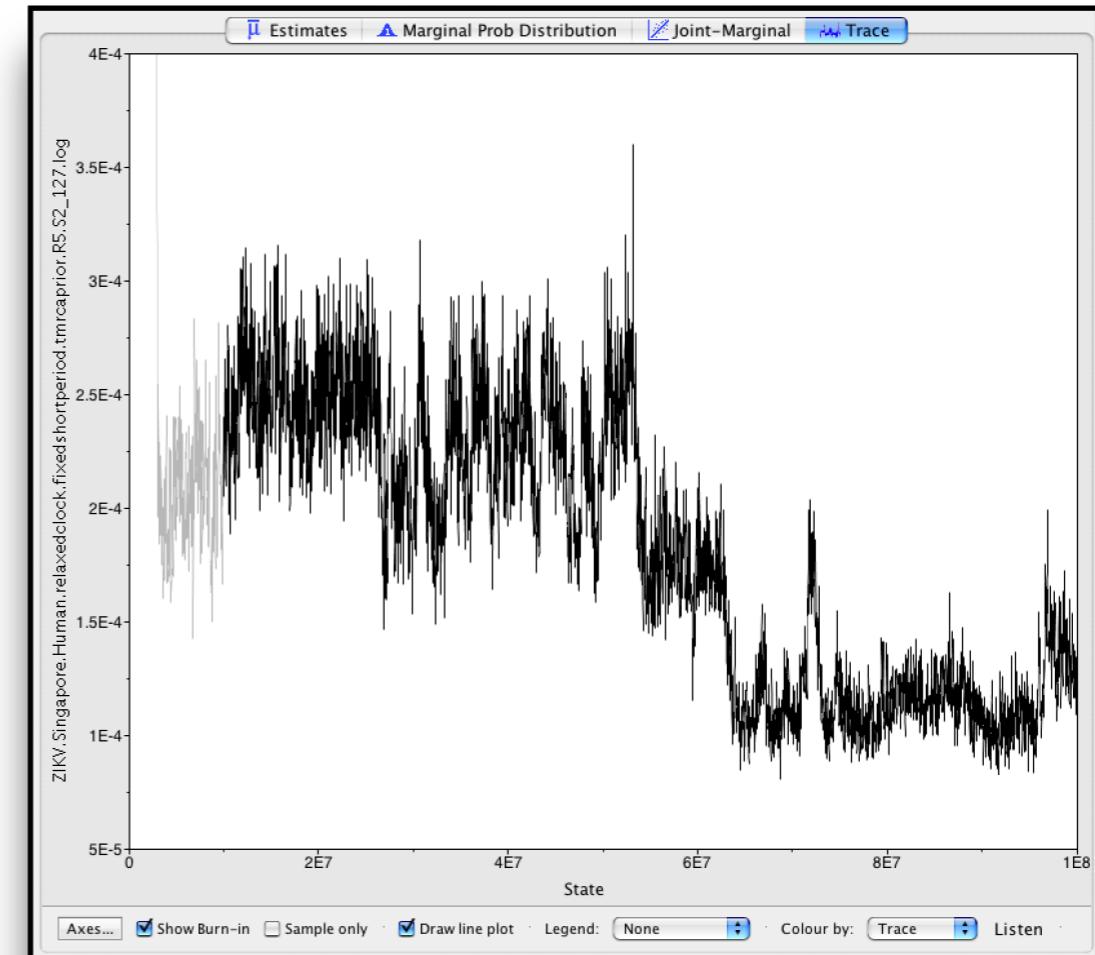


- Demographic reconstructions

# Look at the chains first!



Mixing well! 😊



Not mixing! 😢

- Demographic reconstructions



# TreeAnnotator

(Included with BEAST2)

---

- Analyse trees file from BEAST2 runs
- Produces single summary tree (MCC) with node annotations (including clade posterior probabilities)
- Positions internal nodes according to average taxon set MRCA times in trees file
- Note that the MCC tree is just a heuristic summary: may produce negative branch lengths when topological uncertainty is large!

## Input:

- Tree log file  
(many trees)

## Output:

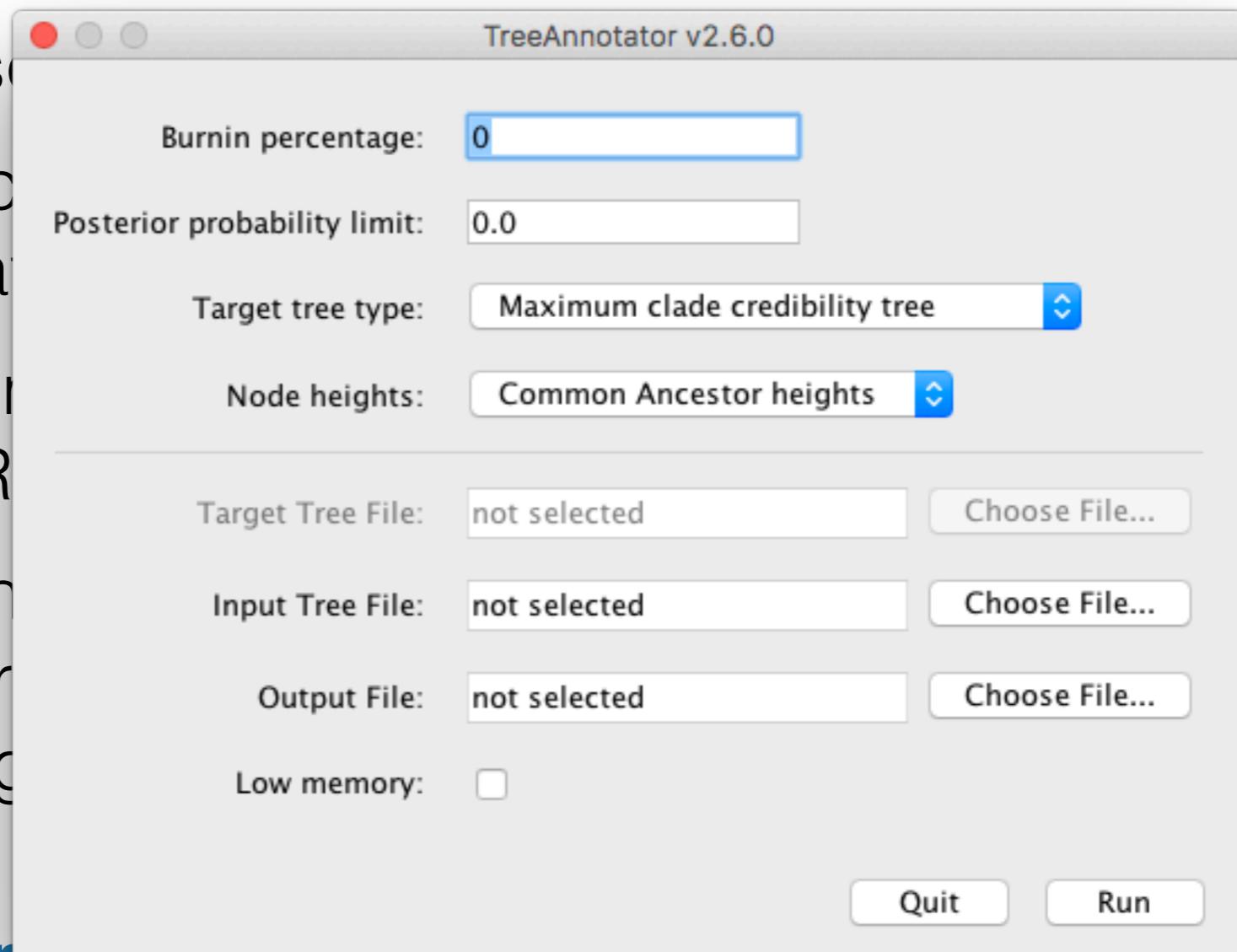
- MCC tree  
(one annotated summary tree)



# TreeAnnotator

(Included with BEAST2)

- Analyses
- Produces annotated trees
- Positions tips on set MCC tree
- Note that may produce topology



## Input.

- Tree log file  
(many trees)

## Output.

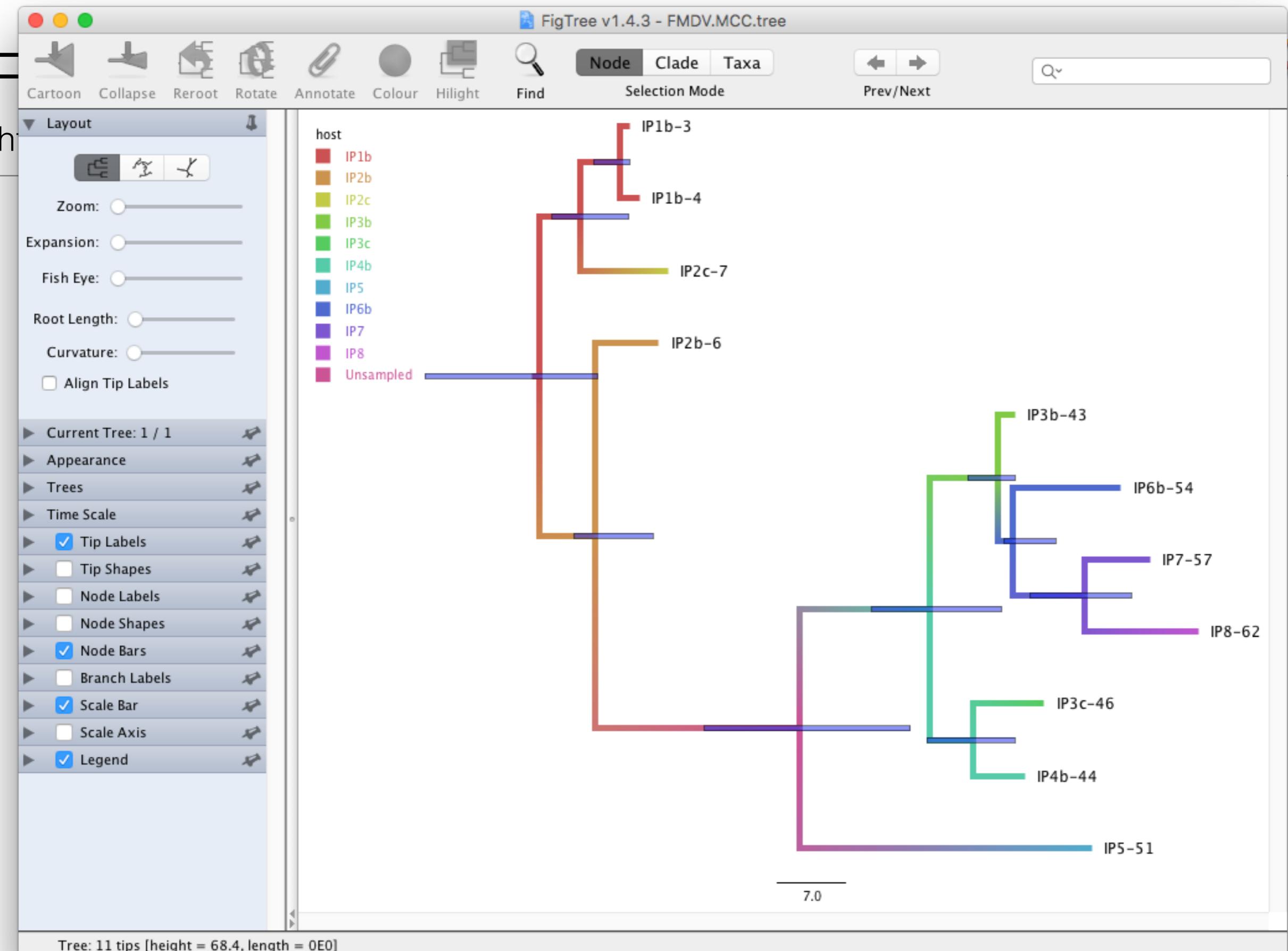
- MCC tree  
(one annotated summary tree)

# FigTree

(<http://tree.bio.ed.ac.uk/software/figtree/>)



- Visualise trees from BEAST2 runs
- Annotate branches and nodes with probabilities and labels



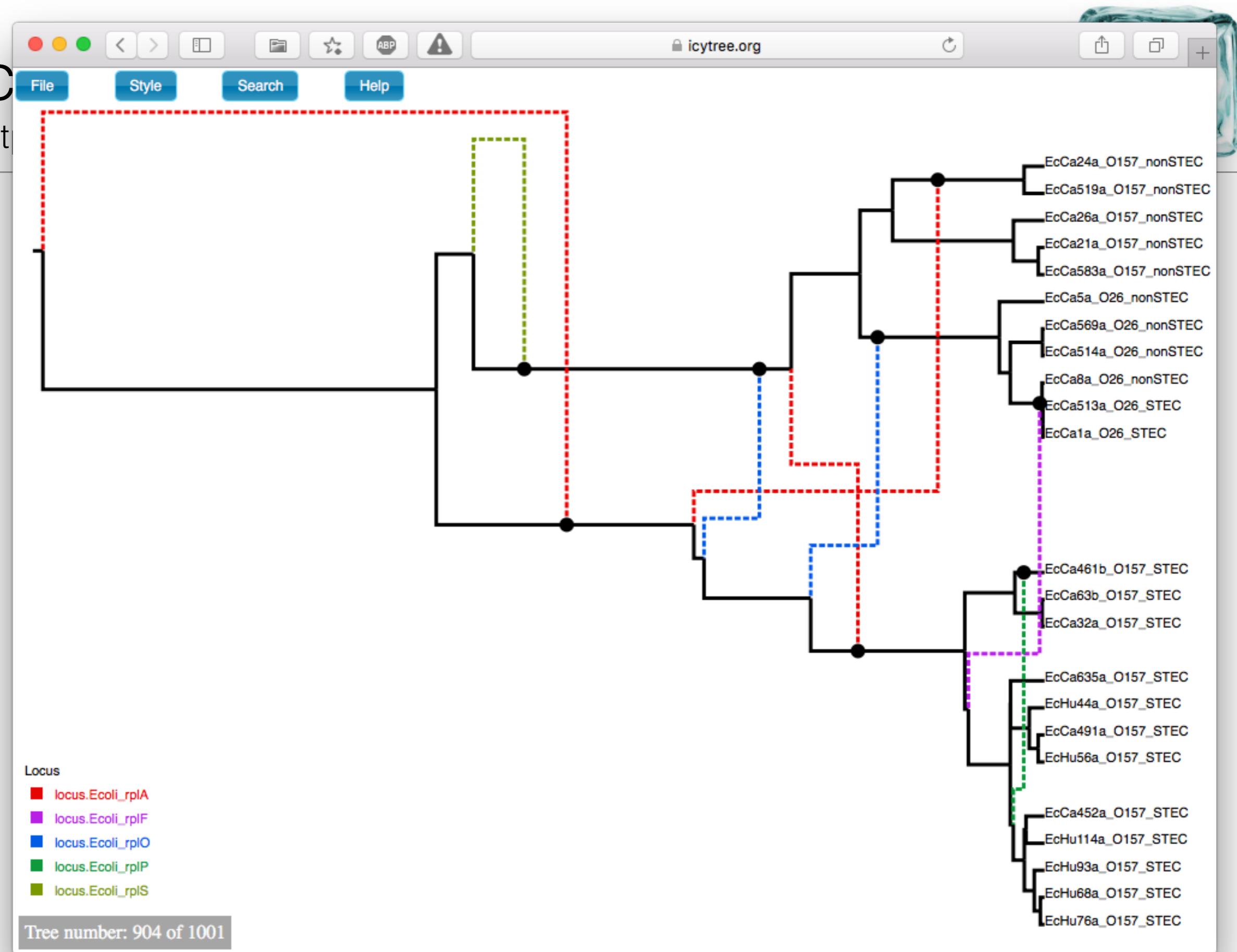
# IcyTree

(<https://icytree.org>)



- Similar to FigTree, but places an emphasis on quick visualisation rather than publication quality output
- Annotate branches and nodes with probabilities and labels
- Better suited for structured models and ancestral recombination graphs (ARGs)
- Faster than FigTree for analysing many trees
- Web app (no installation required)

|C  
(htt

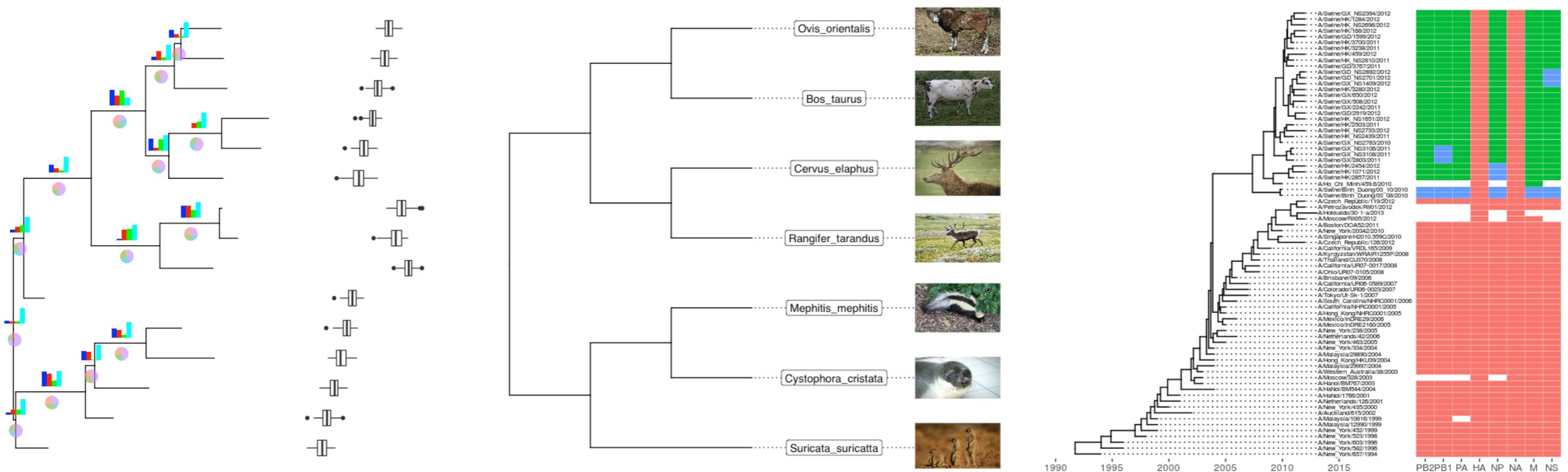


# ggtree

(<https://guangchuangyu.github.io/software/ggtree/>)



- R-package to visualise trees using ggplot grammar
- Works with BEAST2 tree files (and many other packages)
- Can easily annotate trees with other analyses in R



# On the program for today

---

- (1) What is phylodynamics?
- (2) Bayesian inference recap
- (3) BEAST2 introduction

## **Tutorial: Molecular clock dating (part i)**

- (4) Molecular clock models

## **Tutorial: Molecular clock dating (part ii)**

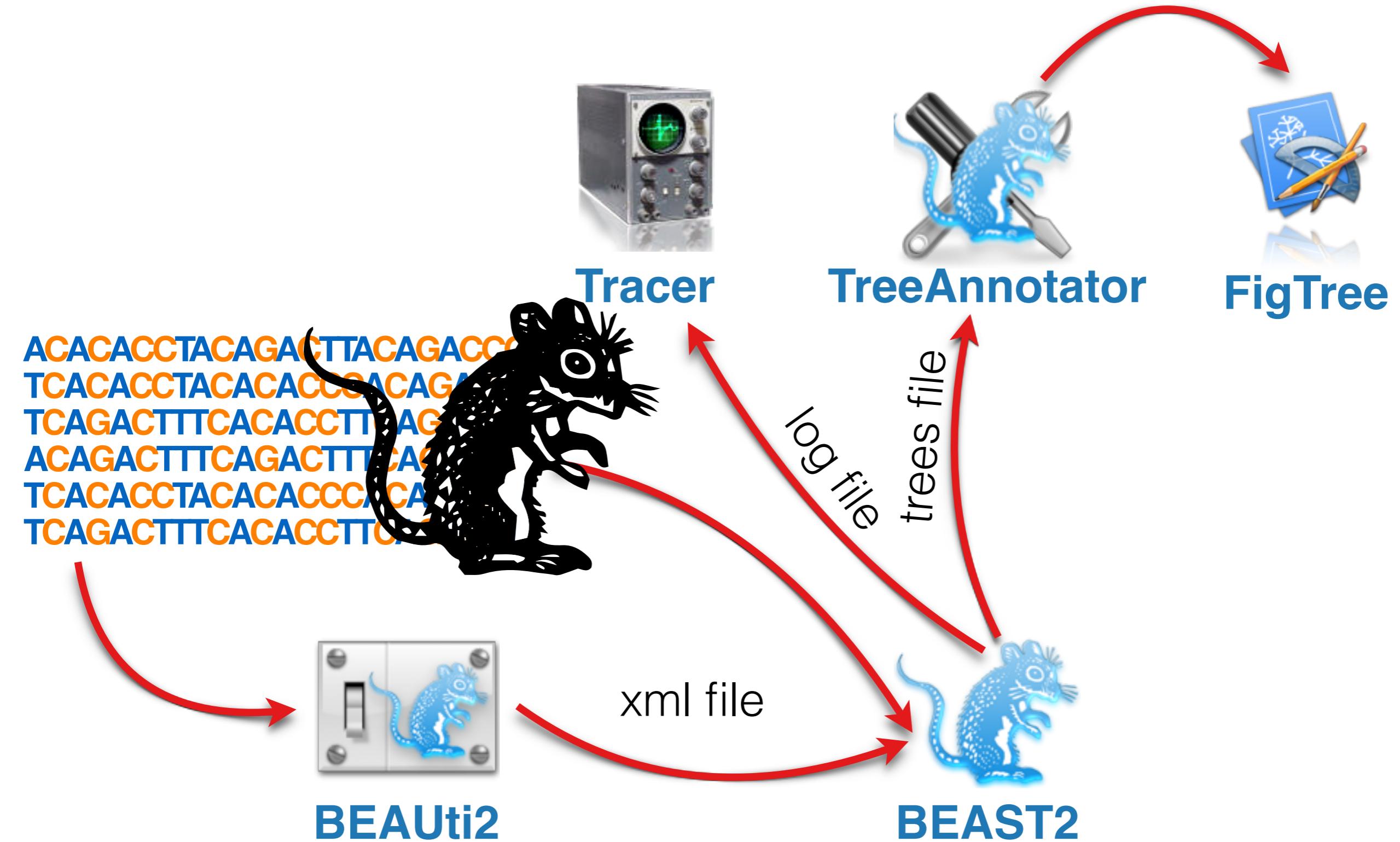
- (5) Setting priors

## **Tutorial: Phylodynamics (part i)**

- (6) Tree priors

## **Tutorial: Phylodynamics (part ii)**

# BEAST2 workflow



# Tools of the trade

---

## **BEAST2**

Software implementing MCMC for model parameter and tree inference

## **BEAUTi2**

Part of BEAST2 package for setting up the input file (.xml)

## **Tracer**

Analysis of BEAST output files (.log)

## **TreeAnnotator**

Analysis of BEAST output files (.trees)

## **FigTree, IcyTree, ggtree**

Visualisation of trees (.trees)

# BEAST best practice

(This is just a guideline and each analysis is unique)

---

## Before you begin

- 1) Know your data  
(check alignment, remove outliers, check clock signal etc.)
- 2) Plan your analysis carefully

## Before you run the analysis

- 3) Ask someone else to look at your XML file
- 4) Sample from the prior (run without data)

## Actually running the analysis

- 5) Run analysis with multiple chains

## After the analysis

- 6) Combine chains
- 7) Assess convergence and mixing
- 8) Ask someone else to look at your log files

# BEAST troubleshooting priors

(Common reasons for lack of joy)

---

**My analysis won't start!** 😞

**P(X = “parameter out of bounds”) = 0.9**

→ Check initial values fall within prior bounds!

**P(X = “numerical underflow”) = 0.05**

→ Increase number of initialisation attempts  
→ Use a different random number seed  
→ **Use a better starting tree**

**P(X = “package not installed”) = 0.02**

→ Install the package...

**P(X = “something weird”) = 0.01**

→ Contact developers

# BEAST troubleshooting priors

(Common reasons for lack of joy)

---

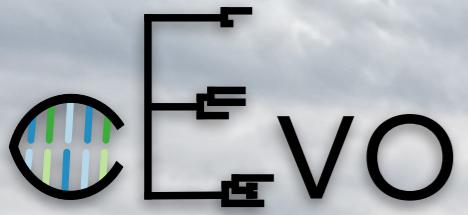
**My analysis won't start!** 😞

**P(X = “output files already exist” | new user) = 0.9**

→ Check “overwrite existing log and tree files”

**P(X = “XML syntax error” | first time editing XML) = 0.9**

- Create template in BEAUti
- Look at examples from the developers
- **Look at the source code**



0111010  
1001100  
0111010  
0111011  
0111010  
0111010

# Molecular clock models

## Louis du Plessis

**ETH** zürich

**DBSSE**

# On the program for today

---

- (1) What is phylodynamics?
- (2) Bayesian inference recap
- (3) BEAST2 introduction

**Tutorial:** Molecular clock dating (part i)

## **(4) Molecular clock models**

**Tutorial:** Molecular clock dating (part ii)

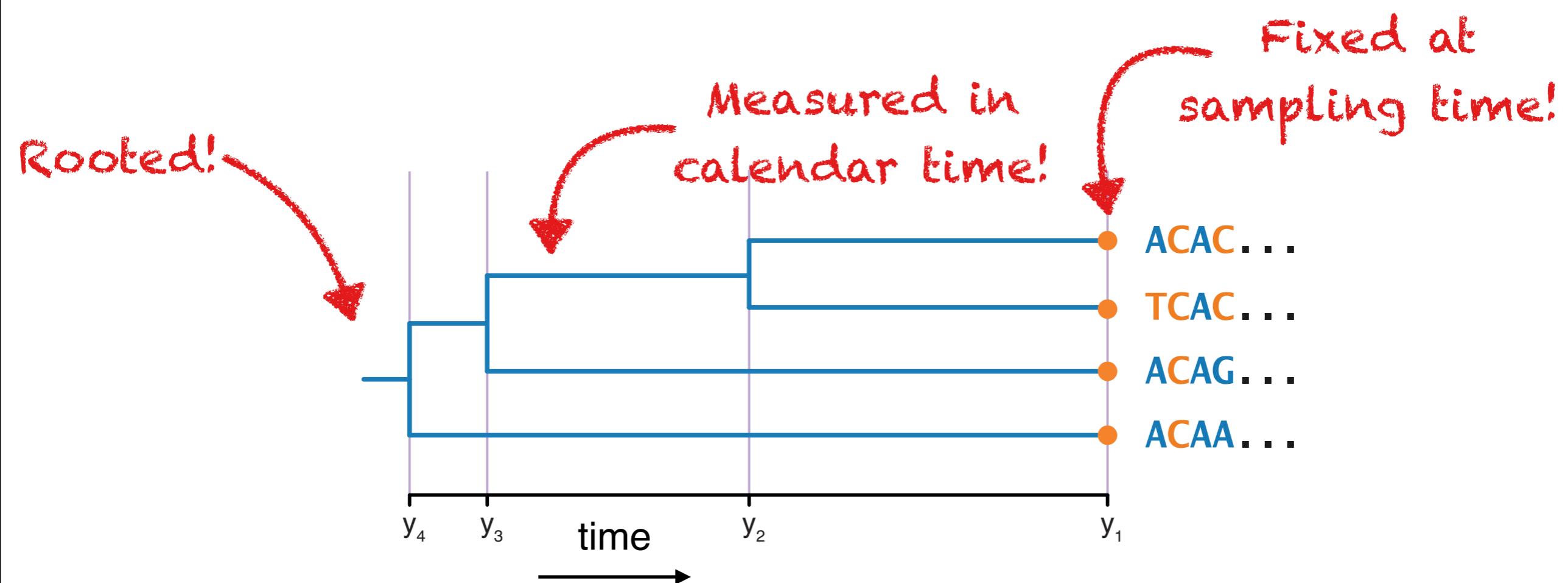
- (5) Setting priors

**Tutorial:** Phylodynamics (part i)

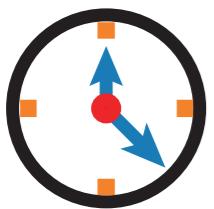
- (6) Tree priors

**Tutorial:** Phylodynamics (part ii)

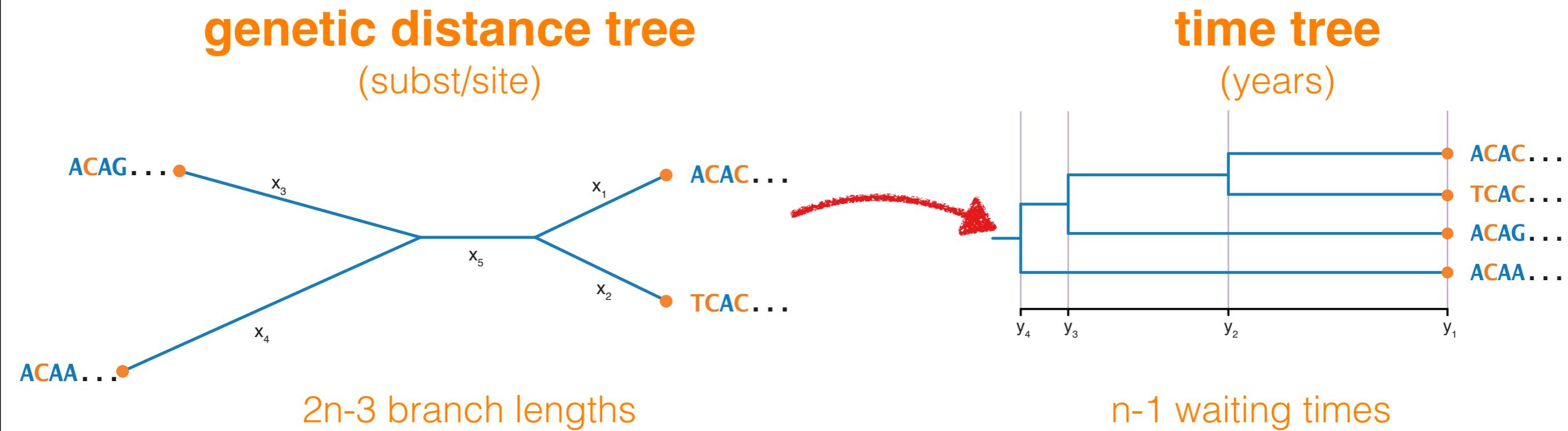
# What is a time tree?



- **Rooted** phylogeny (direction of evolution is known)
- Branch lengths are measured in **calendar time** (units of years, months, days)
- Leaves are fixed at **sampling times**



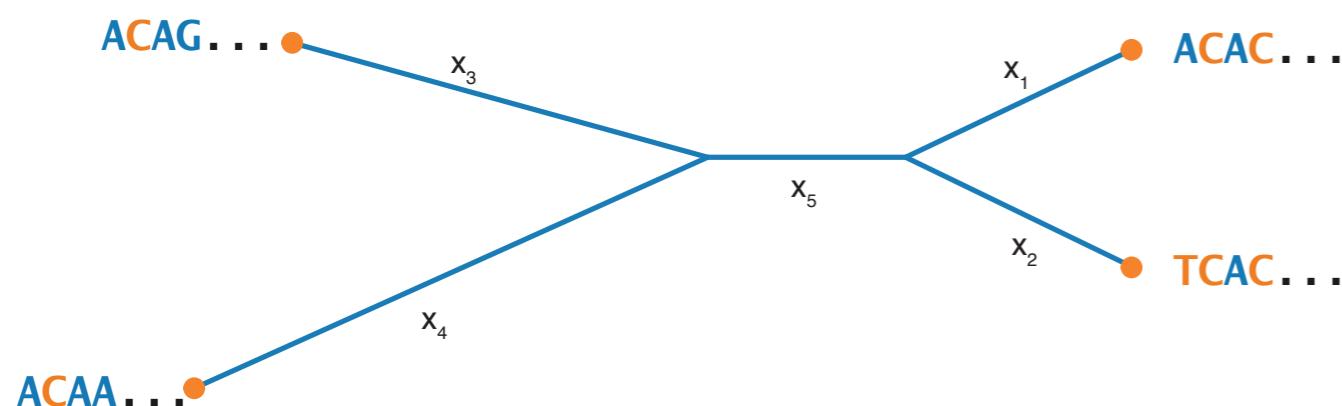
# Molecular clock constraint



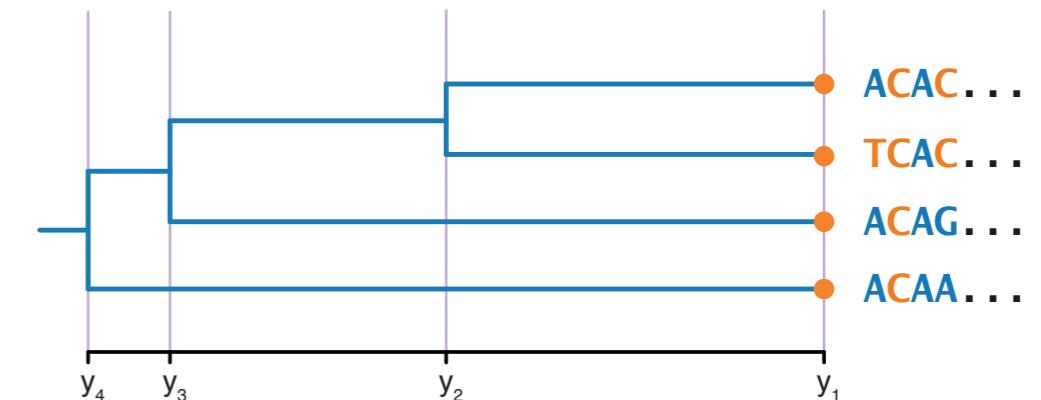
## Molecular clock constraint

- Assume the rate of evolution is constant across all branches
- **Linear** accumulation of substitutions over time

## genetic distance tree

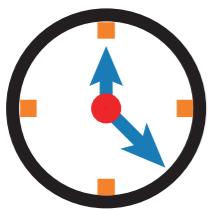


## time tree



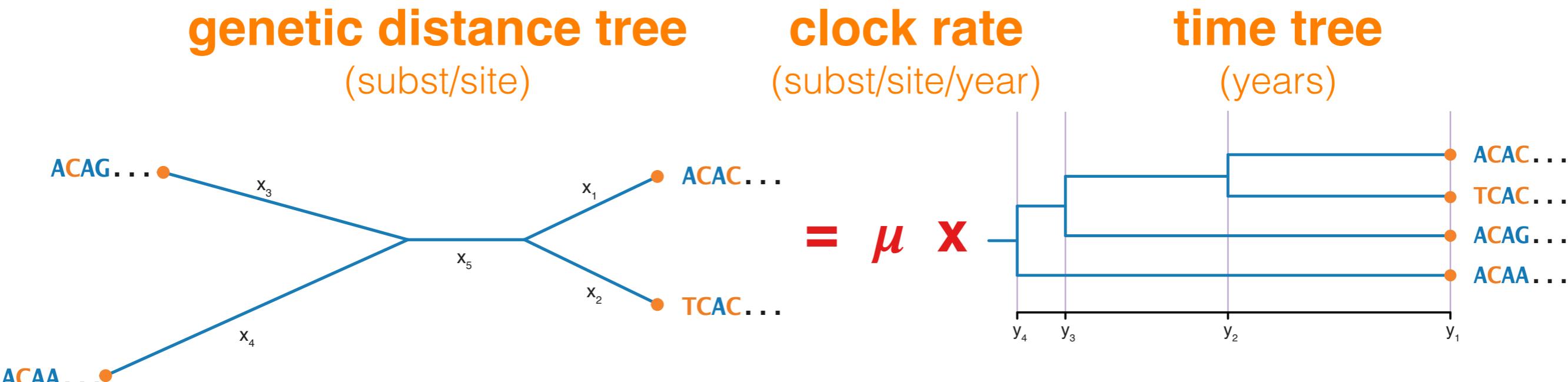
- $2n-3$  **branch lengths**
- Branch lengths are **independent and identically distributed** among branches
- Root location or direction of evolution is not known
- Topology implies nothing about individual branch lengths

- $n-1$  **waiting times**
- Rate of evolution is the **same** on all branches  
(strict clock assumption)
- Root of the tree is equidistant from all tips  
(for contemporaneous samples)
- Topology constrains some branch lengths (2 branches in a cherry must be of equal length)



# Clock rate: $\mu$

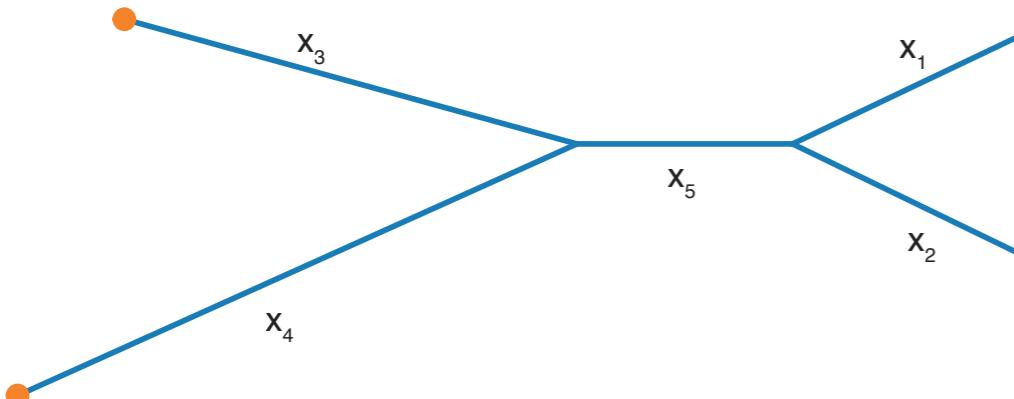
- Rate at which substitutions appear
- Usually measured in substitutions/site/year
- **Genetic distance = Rate x Time**



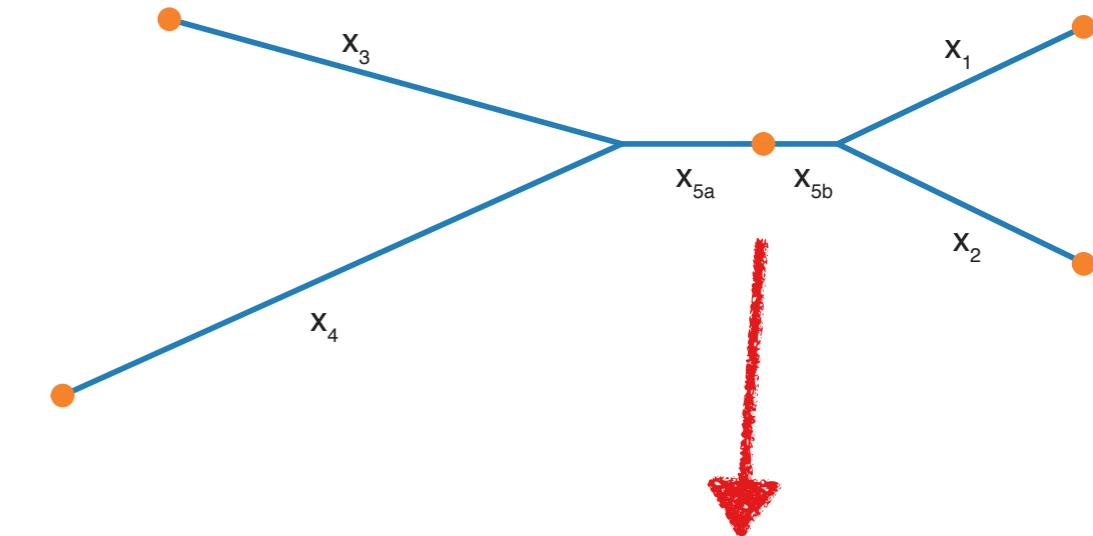
**Clock rate converts branch lengths of time tree into genetic distances**

# Rooting the tree

1) Estimate genetic distance tree

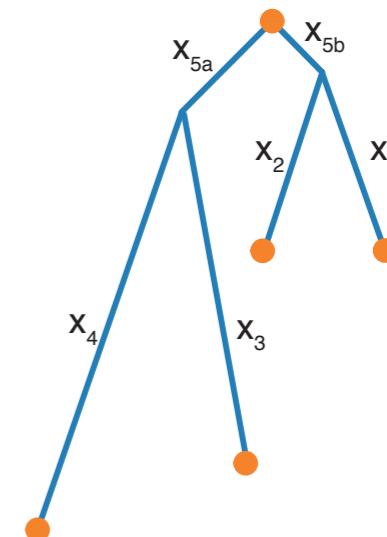


2) Pick root location and split branch

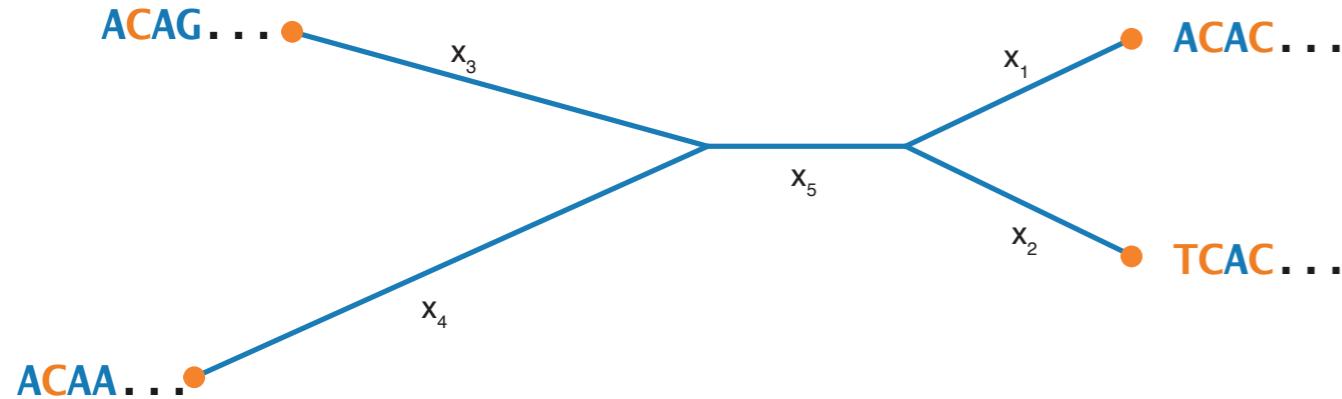


- **The likelihood stays the same!**
  - Substitution model is **reversible**
  - $x_5 = x_{5a} + x_{5b}$
- **Pulley principle:** Root location has no effect on the likelihood

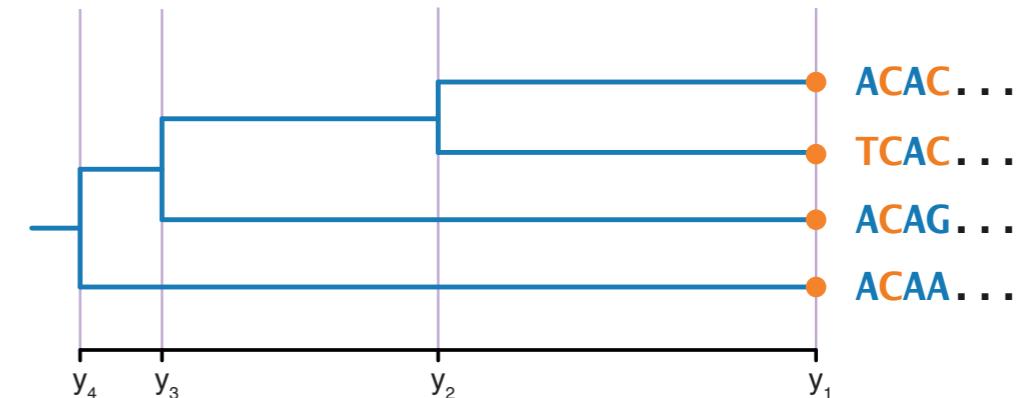
3) PULL!



## genetic distance tree (subst/site)



## time tree (years)



$$P\left(\begin{array}{c} \text{ACAC...} \\ \text{TCAC...} \\ \text{ACAG...} \end{array} \mid \begin{array}{c} \text{tree topology} \\ \text{mutation matrix} \end{array}\right) = P\left(\begin{array}{c} \text{ACAC...} \\ \text{TCAC...} \\ \text{ACAG...} \end{array} \mid \mu \times \begin{array}{c} \text{tree topology} \\ \text{mutation matrix} \end{array}\right)$$

if

$$x_1 = \mu(y_2 - y_1)$$

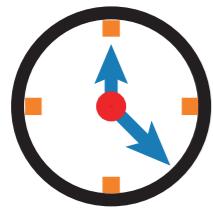
$$x_2 = \mu(y_2 - y_1)$$

$$x_3 = \mu(y_3 - y_1)$$

$$x_4 = \mu[(y_4 - y_1) + (y_4 - y_3)]$$

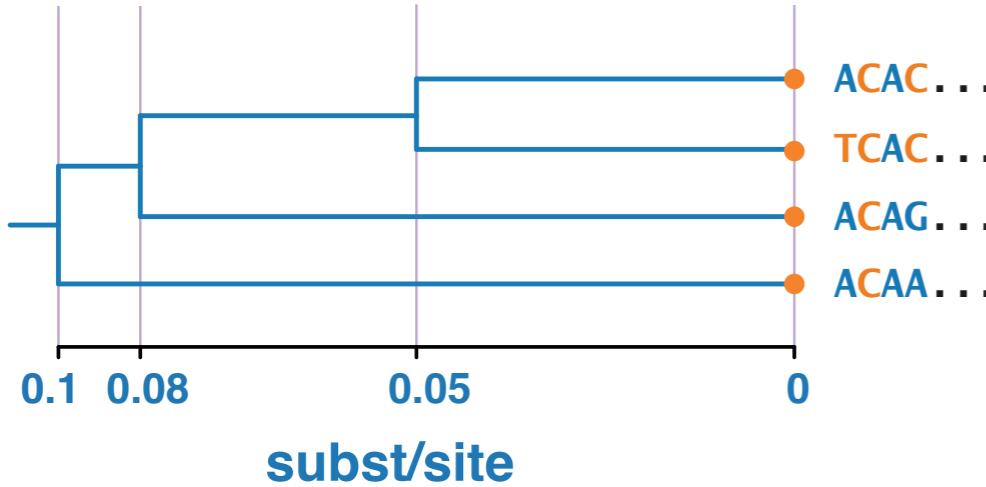
$$x_5 = \mu(y_3 - y_2)$$

**(Pulley principle)**



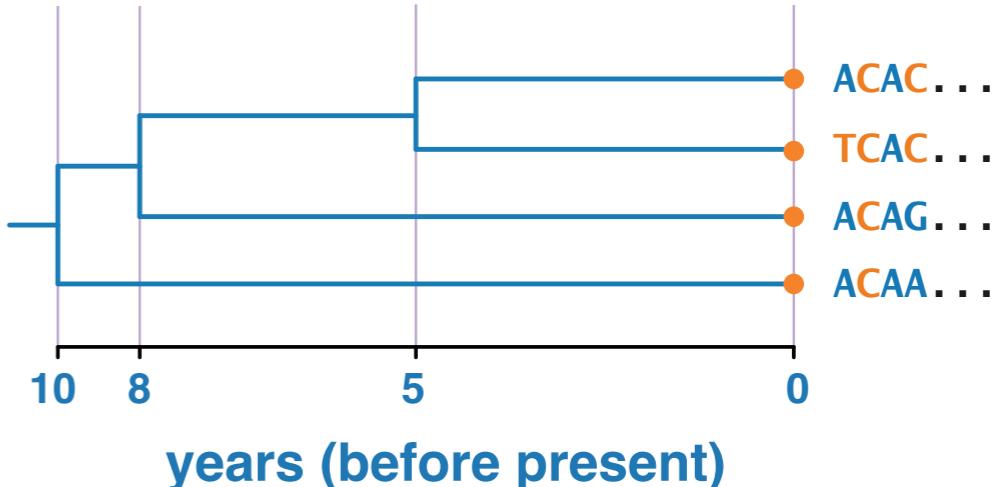
# Non-identifiability of rates and times

**Genetic distance tree**  
(subst/site)



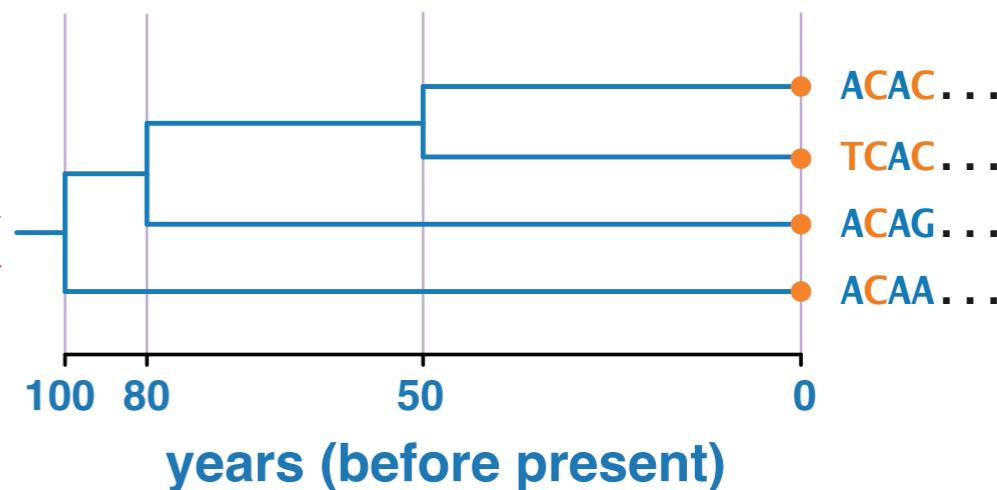
**Clock rate**  
(subst/site/year)

$$= 0.01 \times$$



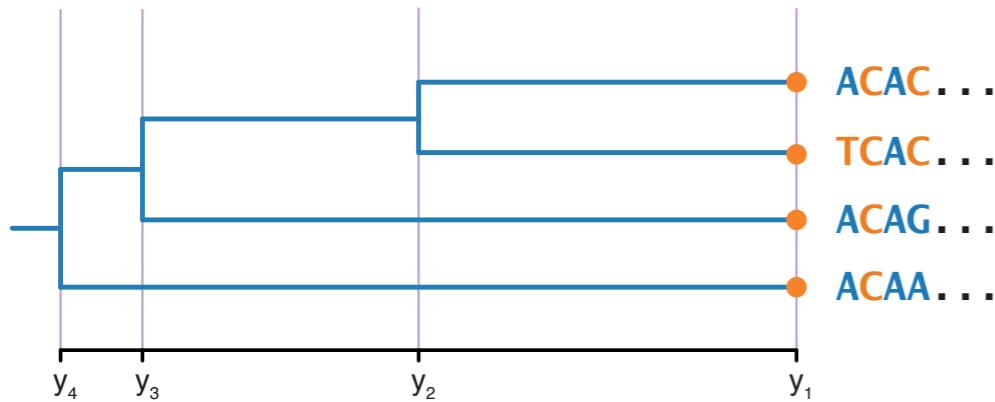
**Time tree**  
(year)

$$= 0.001 \times$$



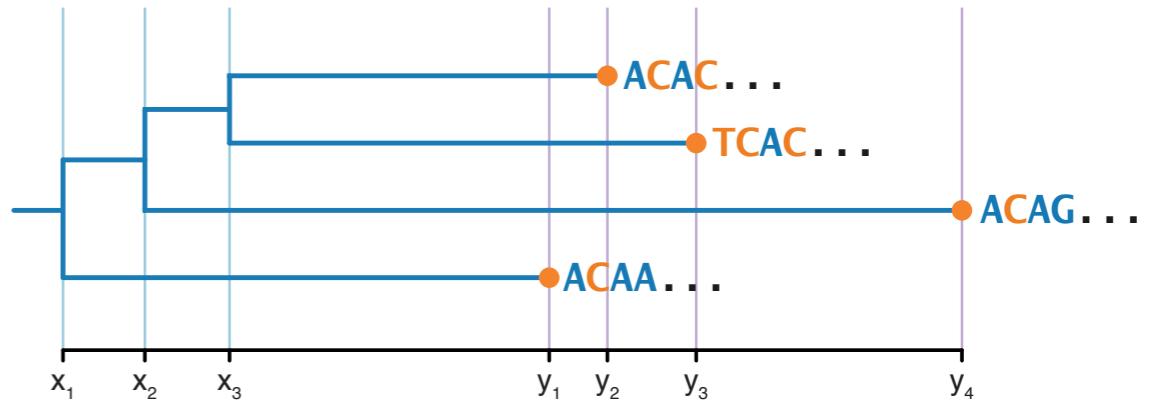
**Need extra information to calibrate the clock!**

## Homochronous data



- All sequences sampled at the same time (contemporaneous)
- Root of the tree is equidistant from all tips
- Non-identifiability of clock rate and time

## Heterochronous data

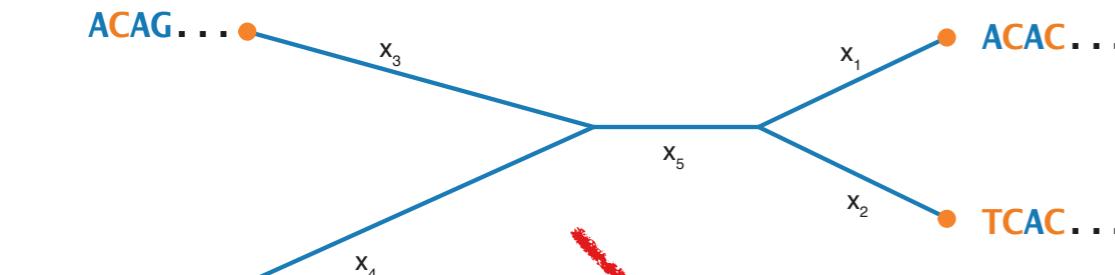


- Sampling period is a significant part of the tree height
- Time to the root is different for each sequence
- **We can use the sampling times to calibrate the clock!**

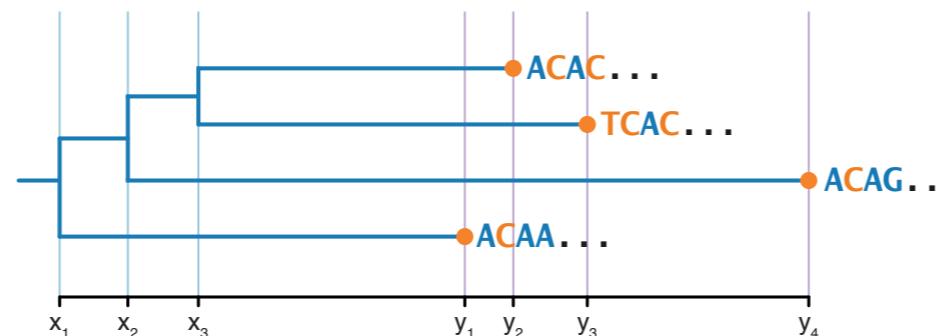
# Calibrating the clock

(heterochronous data)

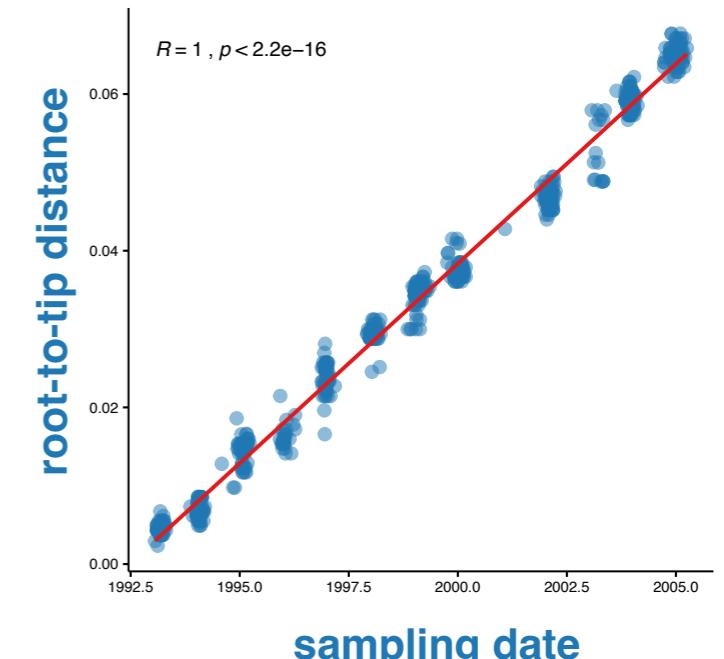
1) Estimate genetic distance tree



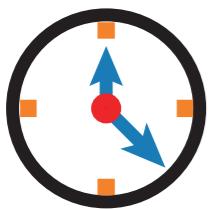
2) Pick a root



3) Plot root-to-tip distance vs sampling date

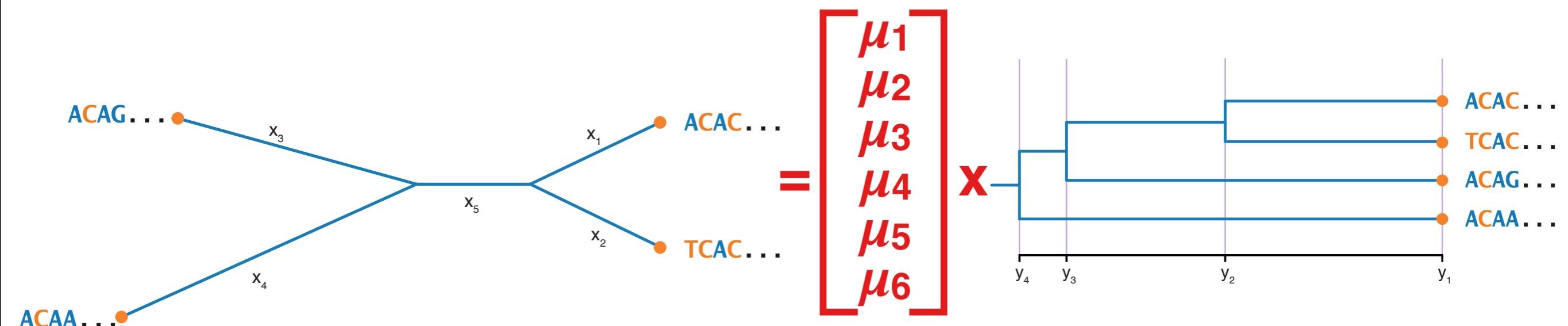


- Use linear regression (least squares) to fit a line through the points ( $\alpha + \beta t$ )
- Clock rate is the slope ( $\beta$ )
- Tree height (tMRCA) is the x-axis intercept ( $-\alpha/\beta$ )
- Repeat steps **2** and **3** to maximise correlation,  $R^2$  etc. or minimise the residuals (only if the true root is not known)



# Relaxed clock: $\bar{\mu}$

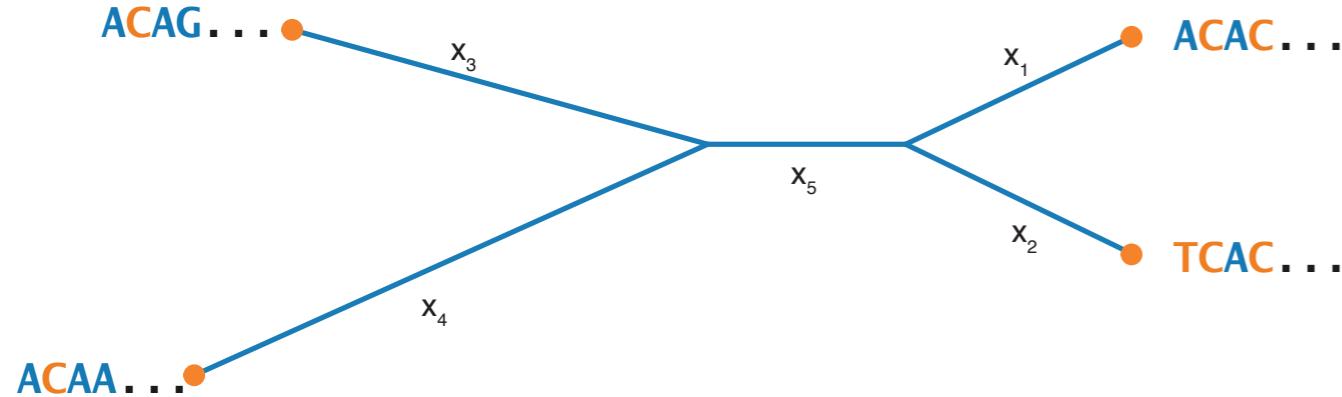
- Different branches can have different clock rates
- **Genetic distance = Rate x Time**



**Clock rate converts branch lengths of time tree into genetic distances**

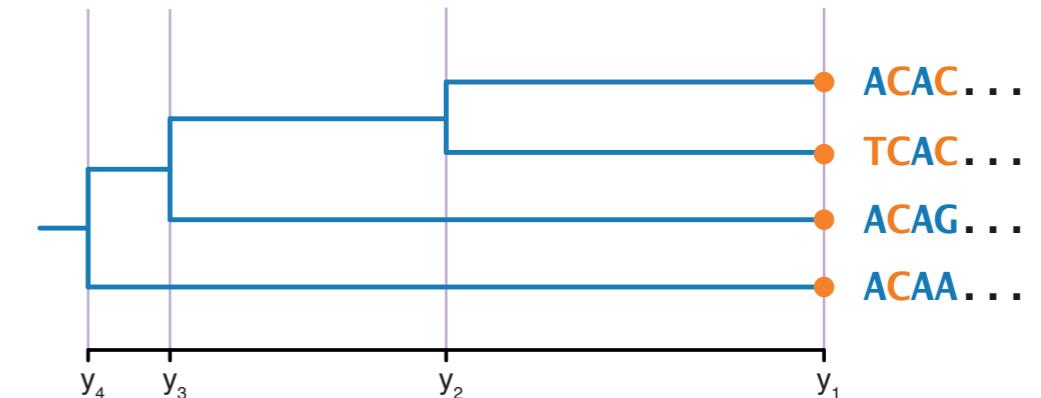
## genetic distance tree

(subst/site)



## time tree

(years)



$$P\left(\begin{array}{c} \text{ACAC....} \\ \text{TCAC....} \\ \text{ACAG....} \end{array} \mid \begin{array}{c} \text{ACAA....} \\ \text{ACAC....} \\ \text{TCAC....} \\ \text{ACAG....} \end{array}, \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{array}, \begin{array}{c} \text{ACAC....} \\ \text{TCAC....} \\ \text{ACAG....} \end{array}\right) = P\left(\begin{array}{c} \text{ACAC....} \\ \text{TCAC....} \\ \text{ACAG....} \end{array} \mid \bar{\mu} \times \begin{array}{c} \text{ACAA....} \\ \text{ACAC....} \\ \text{TCAC....} \\ \text{ACAG....} \end{array}, \begin{array}{c} \text{ACAA....} \\ \text{ACAC....} \\ \text{TCAC....} \\ \text{ACAG....} \end{array}\right)$$

if

$$x_1 = \bar{\mu}_1(y_2 - y_1)$$

$$x_2 = \bar{\mu}_2(y_2 - y_1)$$

$$x_3 = \bar{\mu}_3(y_3 - y_1)$$

$$x_4 = \bar{\mu}_4(y_4 - y_1) + \bar{\mu}_6(y_4 - y_3)$$

$$x_5 = \bar{\mu}_5(y_3 - y_2)$$

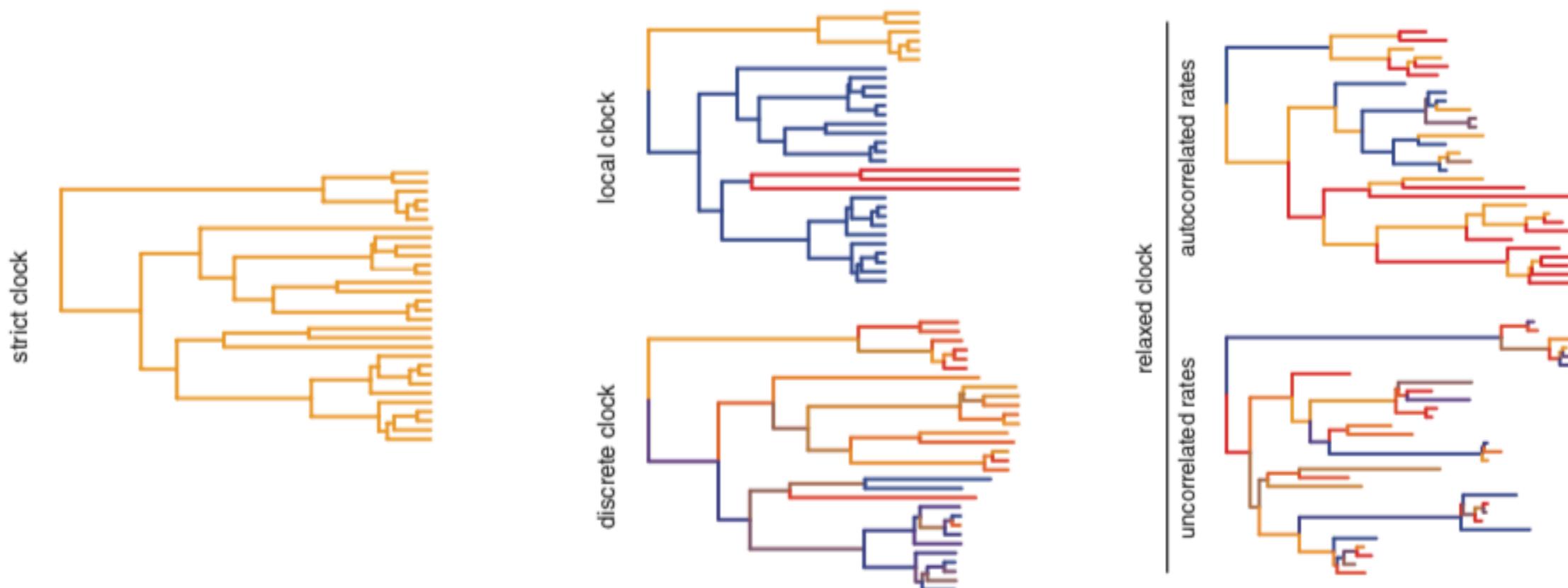
(Pulley principle)

# Causes of rate variation

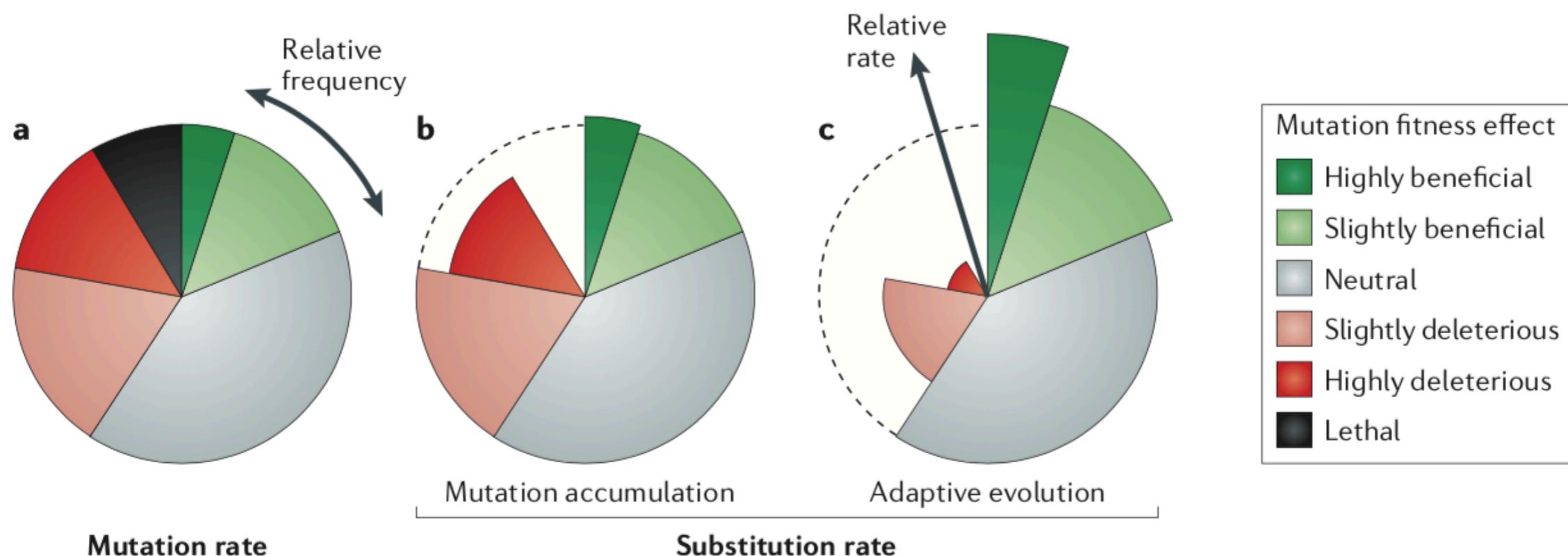
## Differences in

- Generation time
- Metabolic rate
- Error-rate and repair mechanisms
- Population size
- Selective pressure

} **Lineage effects**



# What are we actually estimating?



# On the program for today

---

- (1) What is phylodynamics?
- (2) Bayesian inference recap
- (3) BEAST2 introduction

**Tutorial:** Molecular clock dating (part i)

- (4) Molecular clock models

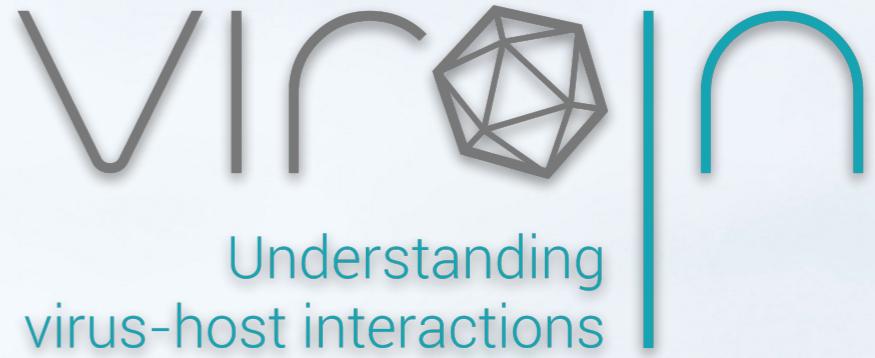
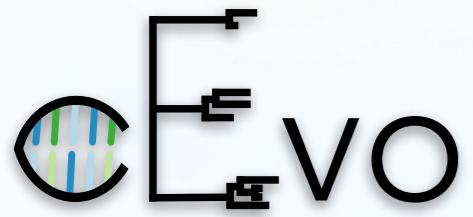
**Tutorial: Molecular clock dating (part ii)**

- (5) Setting priors

**Tutorial:** Phylodynamics (part i)

- (6) Tree priors

**Tutorial:** Phylodynamics (part ii)



0111010  
0100110  
1001010  
0111010  
0111011



Setting priors

Louis du Plessis

ETH zürich

DBSSE

# On the program for today

---

- (1) What is phylodynamics?
- (2) Bayesian inference recap
- (3) BEAST2 introduction

**Tutorial:** Molecular clock dating (part i)

- (4) Molecular clock models

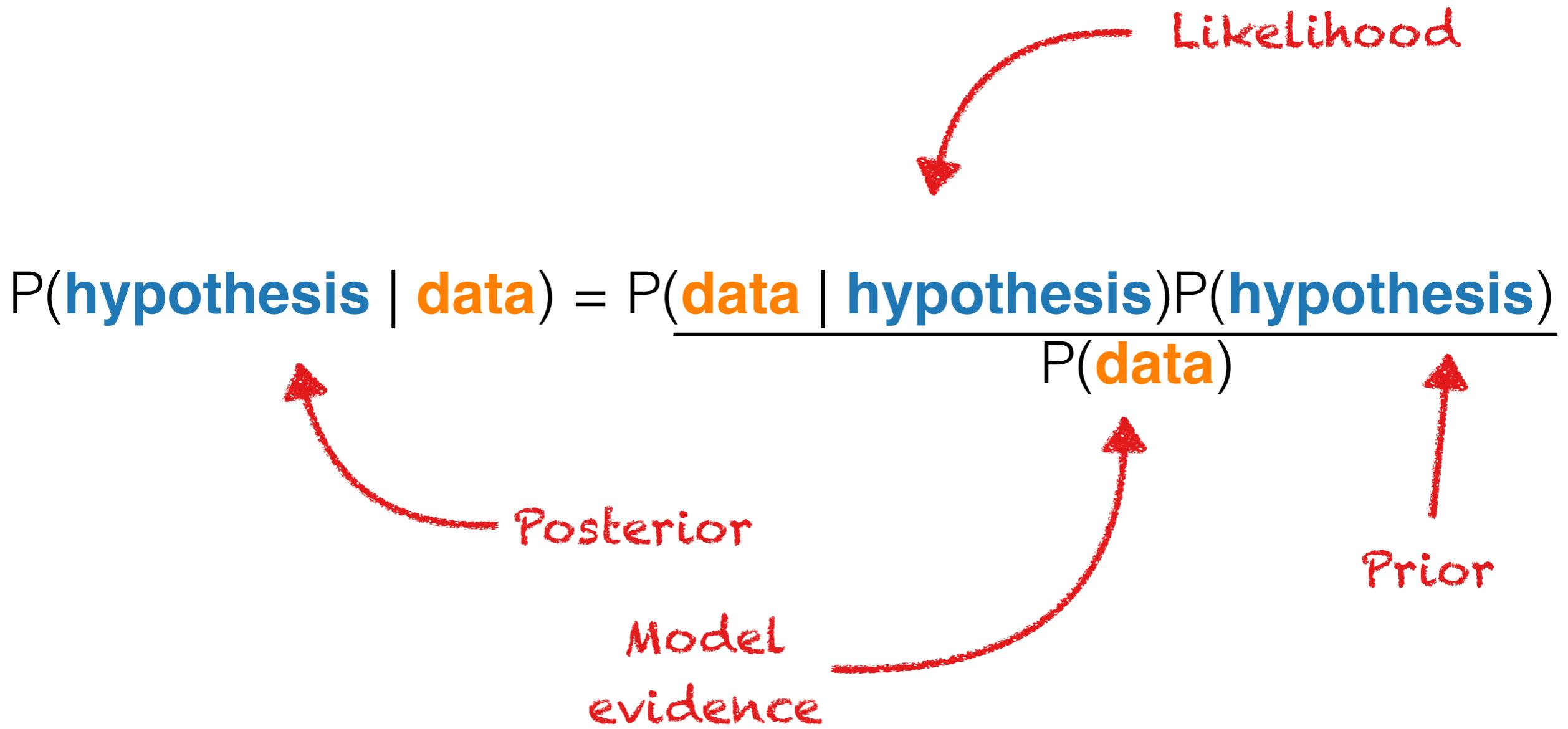
**Tutorial:** Molecular clock dating (part ii)

## **(5) Setting priors**

**Tutorial:** Phylodynamics (part i)

- (6) Tree priors

**Tutorial:** Phylodynamics (part ii)



### Prior $\rightarrow P(\text{hypothesis})$

- Original probability for the model parameters/components (before data collected or analysed)
- Belief in our hypothesis
- All parameters have priors, whether you specify them or not!

# What is a **prior** probability?

---

**The probability of whatever we are interested in, in the absence of possibly relevant data**

- Usually define a **prior distributions** for model **parameters**  
e.g. clock rate, tree, ...
- Often use **parametric distributions**  
e.g. Uniform, Normal, Gamma, Beta, Lognormal, Laplace, ...
- Sometimes a **prior** on a model component  
e.g. substitution model (HKY, GTR, JC, ...)
- Priors can have priors which can have priors *ad infinitum*  
**(hyperpriors)**
- Parameter bounds are part of the prior  
e.g. normal distribution with lower bound 0.

# Improper priors

---

Suppose we want to pick a prior for the molecular clock rate,  $\mu$ , and we want to use an uninformative prior

- What about a uniform prior (constant)?  $f(\mu) = c$

## **Is this uninformative?**

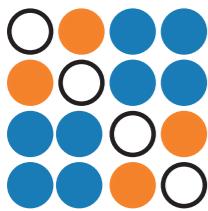
- But this is defined from 0 to infinity!
- Places almost all probability on very large (and biologically unrealistic) values!
- There is no normalising constant such that it integrates to 1!
- $f(\mu) = 1/\mu$  is a better choice for rates (uniform in log space)
  - But it's also impossible to find a normalising constant!
- Probably better to use something like a lognormal distribution
  - Always  $>0$ , long tail, integrates to 1

# How to pick a **prior**?

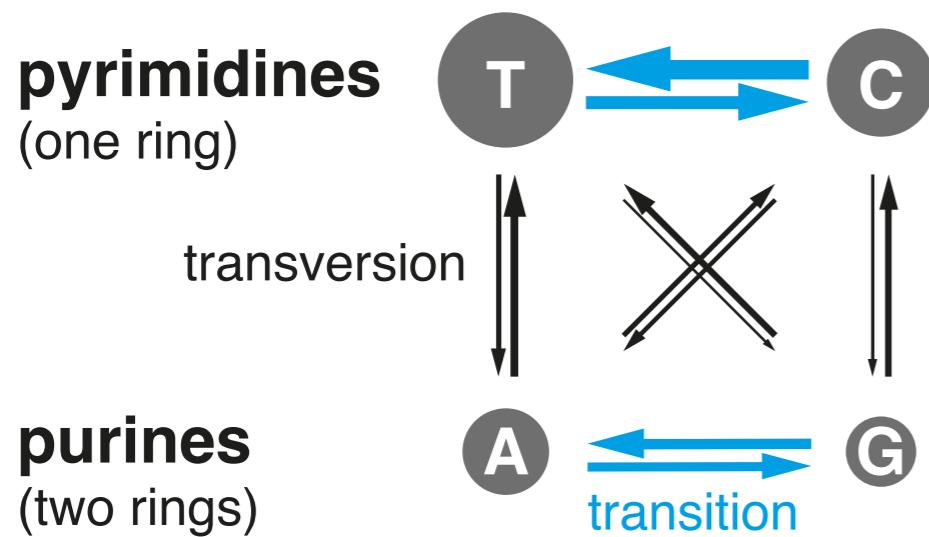
---

- Should be chosen based on your beliefs about model parameters (from independent evidence/experiments)
- Should be chosen with a particular analysis in mind (no priors are universal!)
- Think carefully about plausible ranges of parameters
- Be conservative if you are unsure about the parameter (use diffuse priors)
- Conjugate priors?  
(no — those are just for mathematical convenience)
- Reference priors?  
(yes — if you are an objective Bayesian)
- Do not use improper priors if you can help it  
(priors that integrate to infinity, e.g.  $1/x$ )





# Example: HKY-model (HKY85)

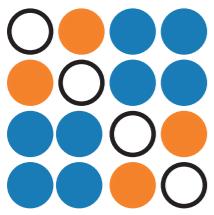


$$(\pi_T, \pi_C, \pi_A, \pi_G)$$

$$\begin{matrix}
 & \text{T} & \text{C} & \text{A} & \text{G} \\
 \text{T} & \cdot & \alpha\pi_C & \beta\pi_A & \beta\pi_G \\
 \text{C} & \alpha\pi_T & \cdot & \beta\pi_A & \beta\pi_G \\
 \text{A} & \beta\pi_T & \beta\pi_C & \cdot & \alpha\pi_G \\
 \text{G} & \beta\pi_T & \beta\pi_C & \alpha\pi_A & \cdot
 \end{matrix}$$

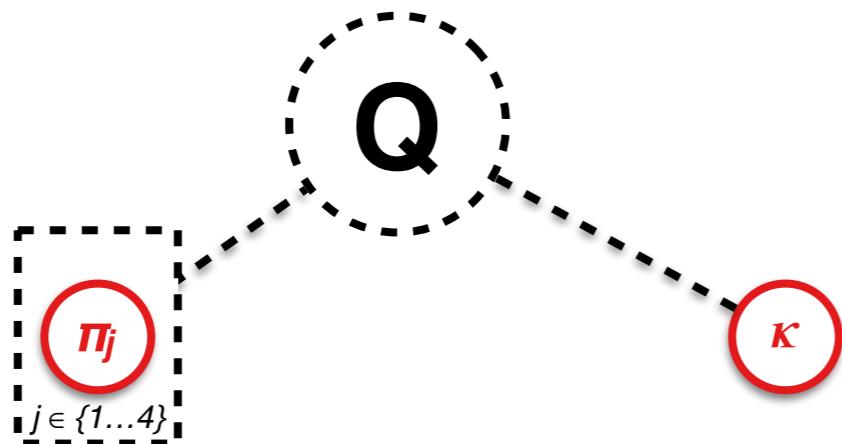
$$= \begin{pmatrix} \cdot & \alpha & \beta & \beta \\ \alpha & \cdot & \beta & \beta \\ \beta & \beta & \cdot & \alpha \\ \beta & \beta & \alpha & \cdot \end{pmatrix} \cdot \begin{pmatrix} \pi_T & 0 & 0 & 0 \\ 0 & \pi_C & 0 & 0 \\ 0 & 0 & \pi_A & 0 \\ 0 & 0 & 0 & \pi_G \end{pmatrix}$$

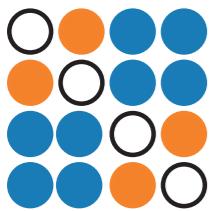
- **Q** matrix gives relative rates of substitution between nucleotides
- 5 parameters:
  - $\kappa = \alpha/\beta$
  - $\pi_T, \pi_C, \pi_A, \pi_G$  — assume these all have the same prior probability



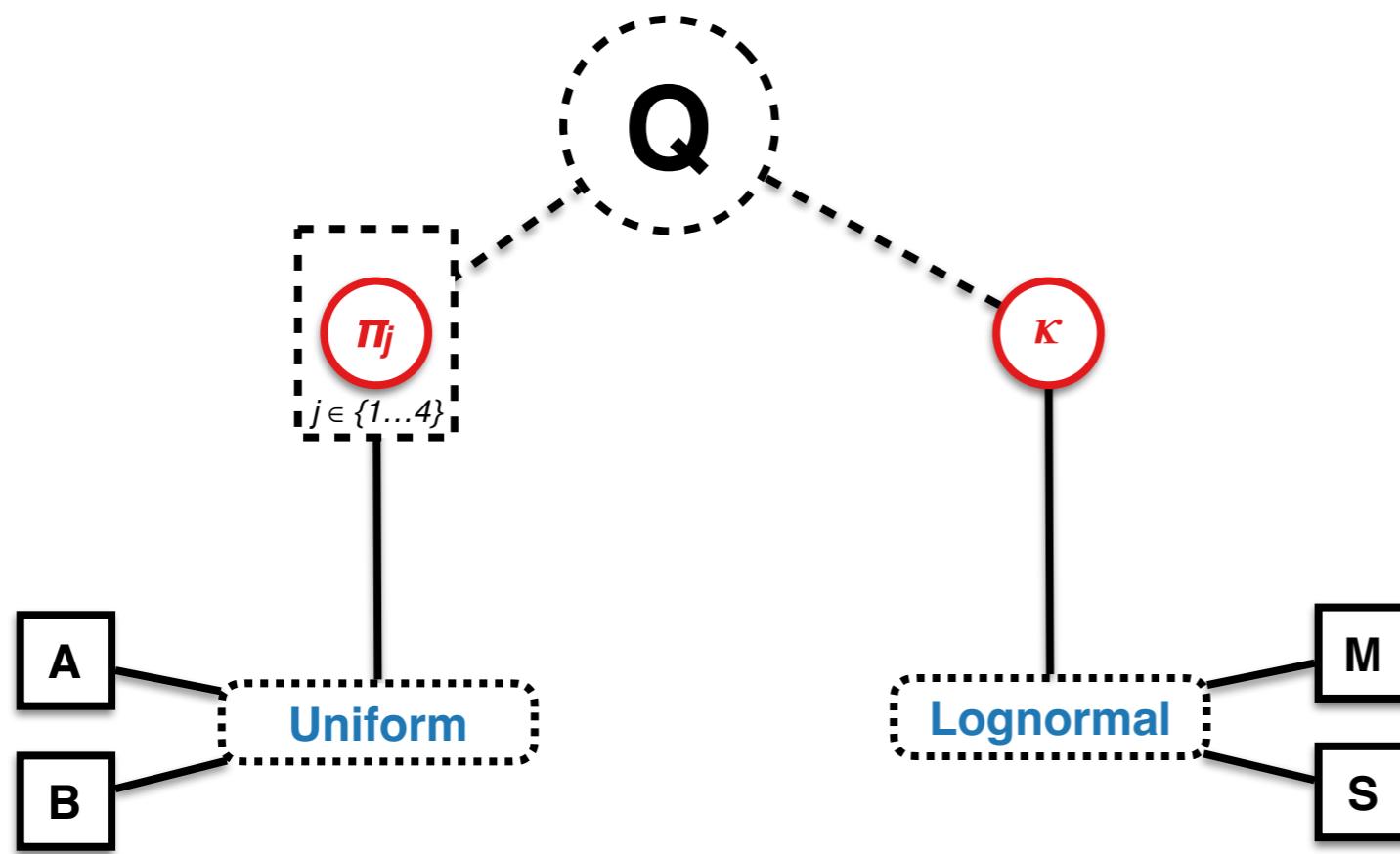
# Example: HKY-model (HKY85)

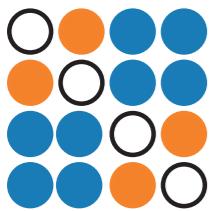
---



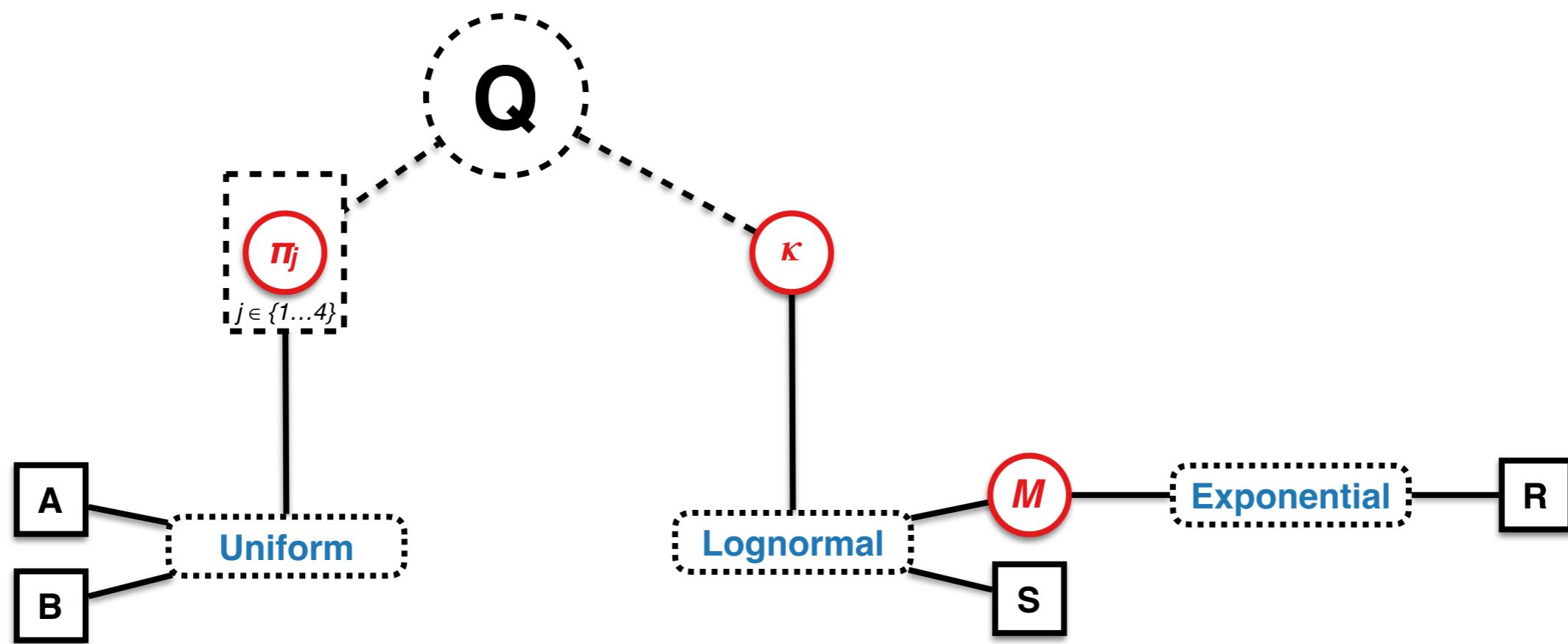


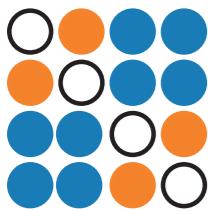
# Example: HKY-model (HKY85)



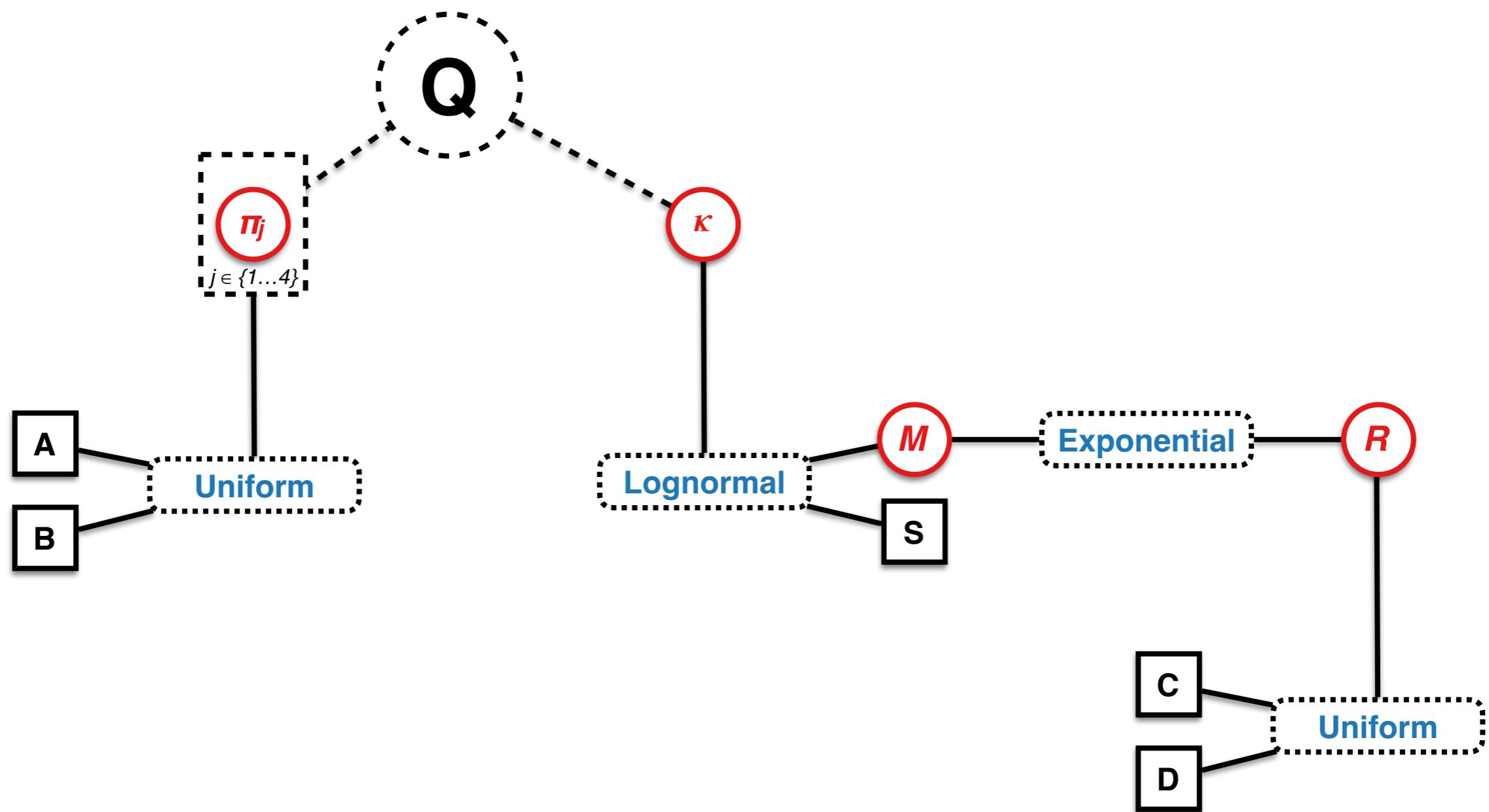


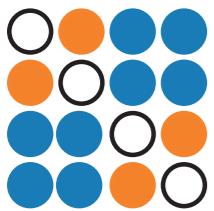
# Example: HKY-model (HKY85)



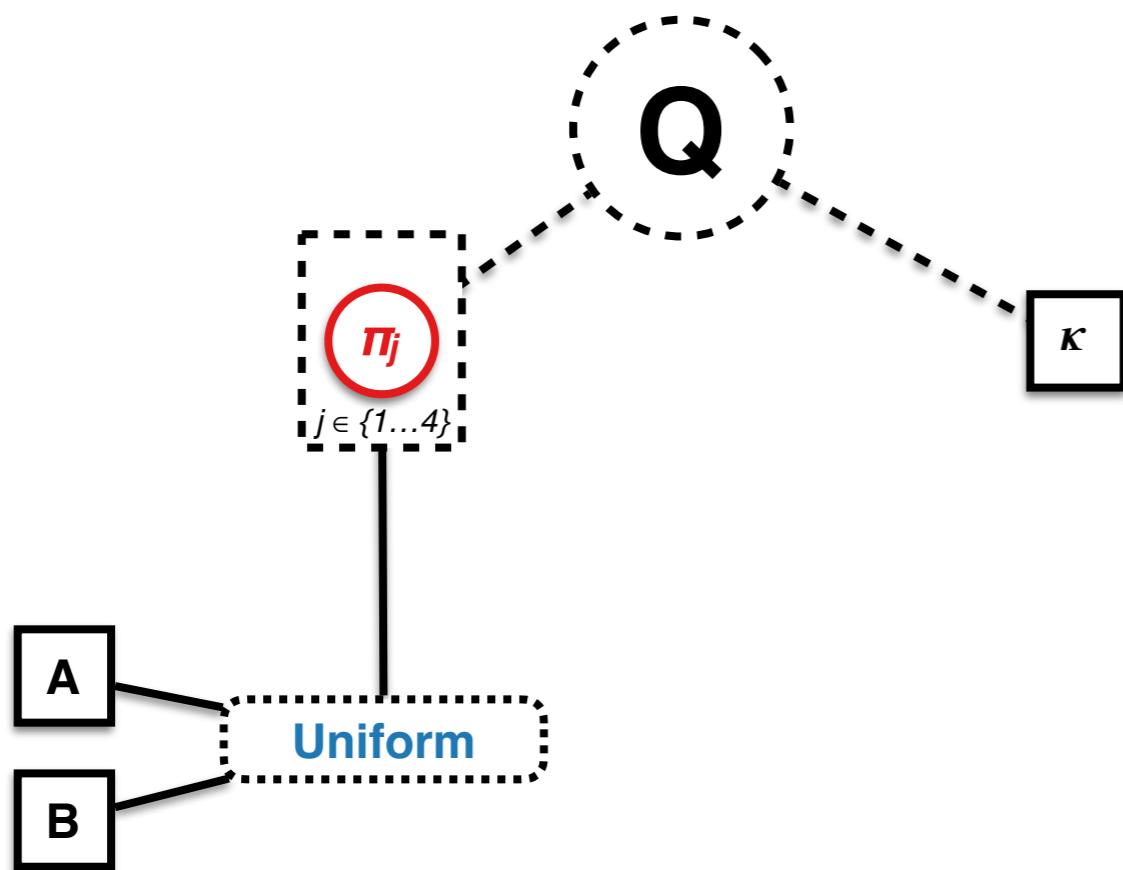


# Example: HKY-model (HKY85)

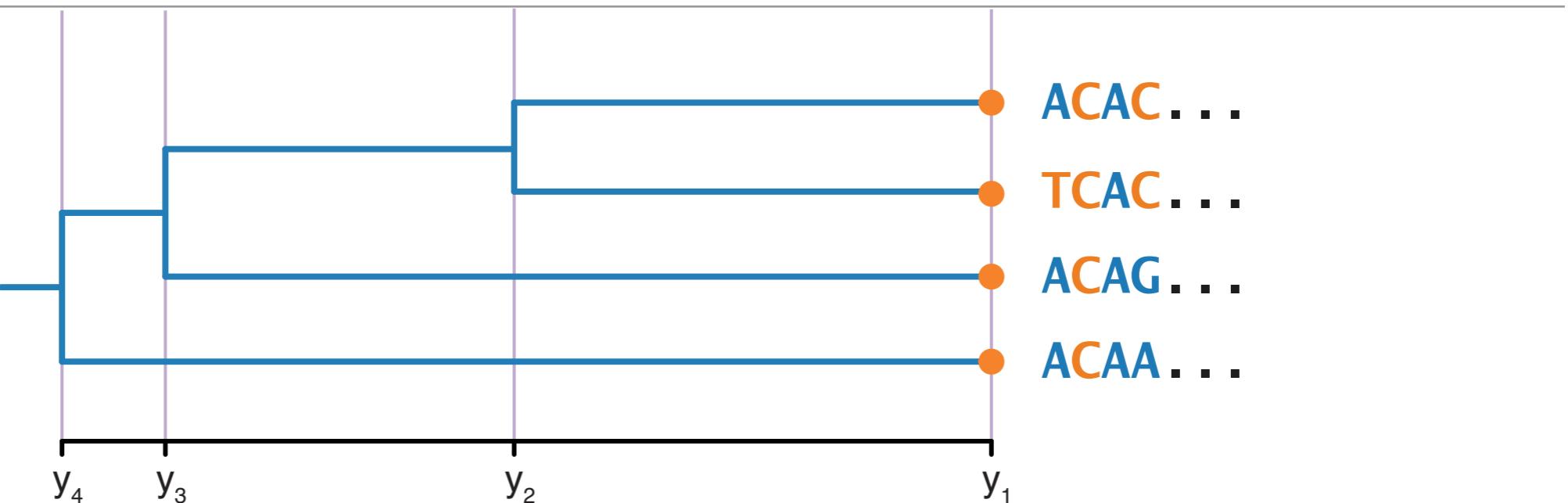




# Example: HKY-model (HKY85)



# Calibration nodes

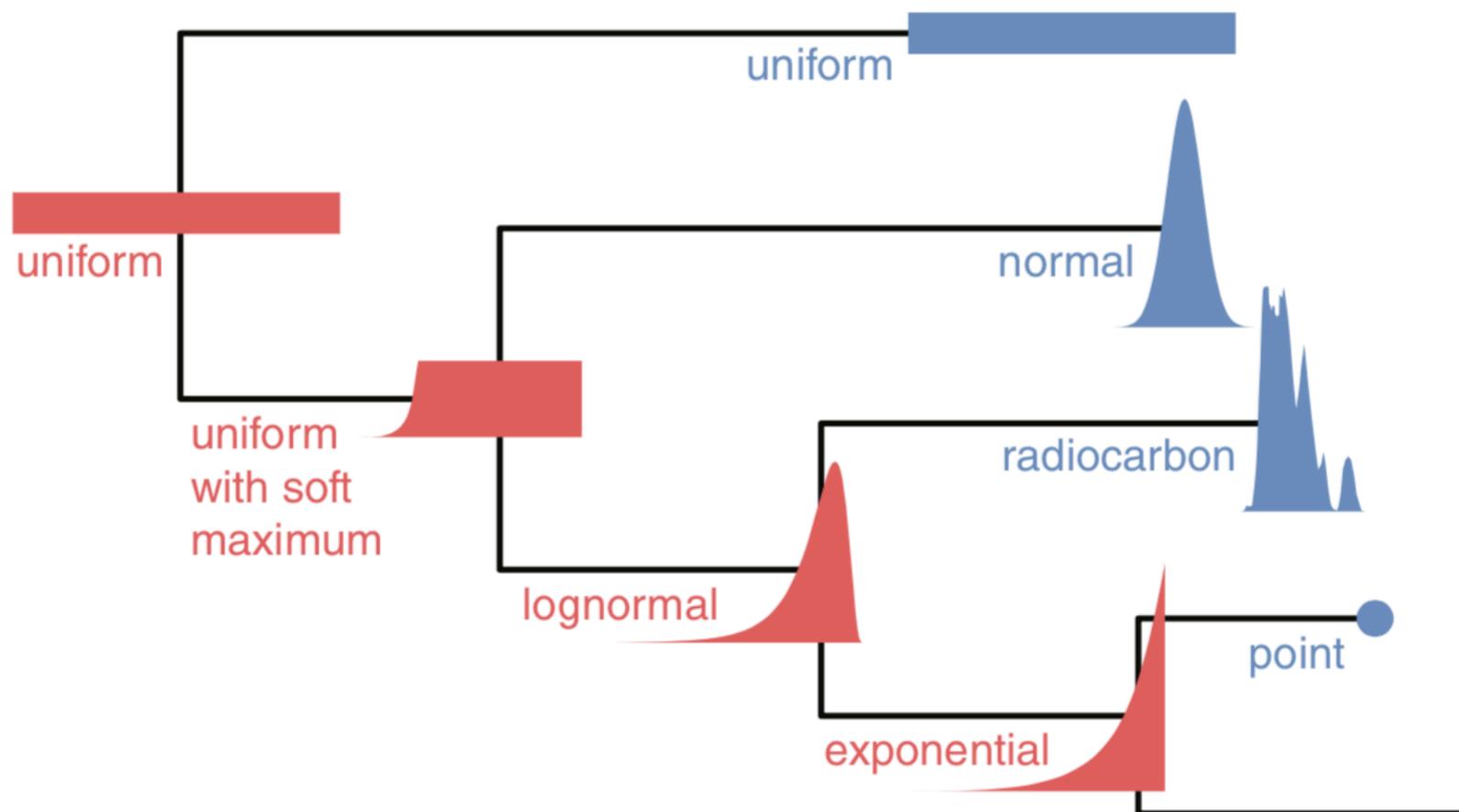


Homochronous trees need external (**prior**) information to date:

- Fix the clock rate
- Use a **calibration node** (or nodes)

# Calibration nodes

---

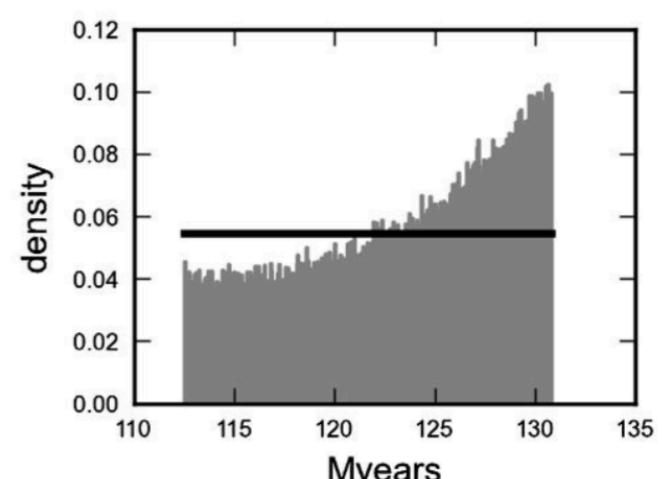
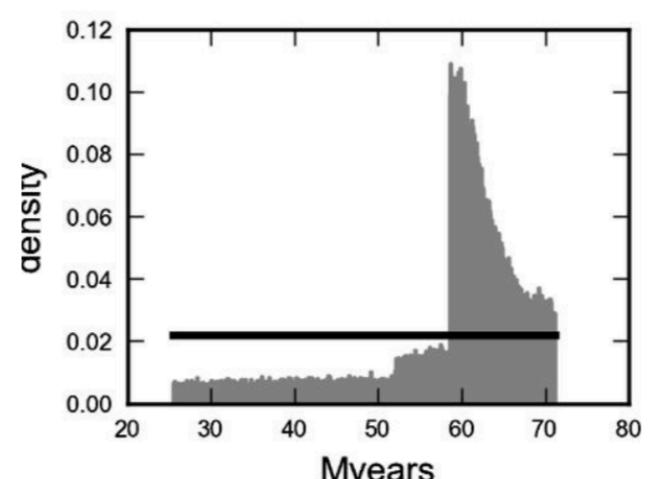
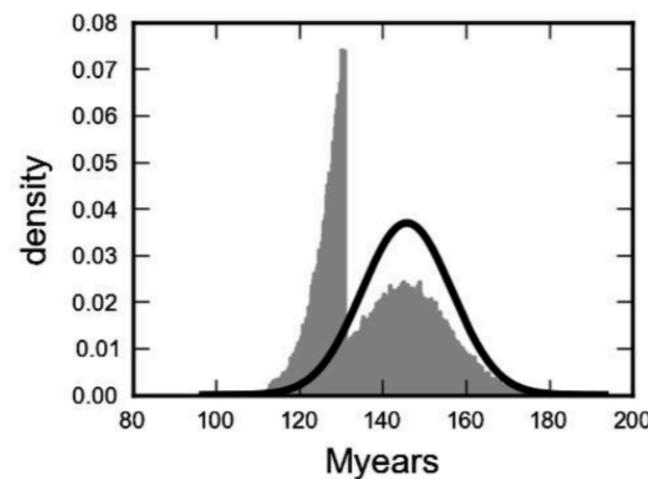
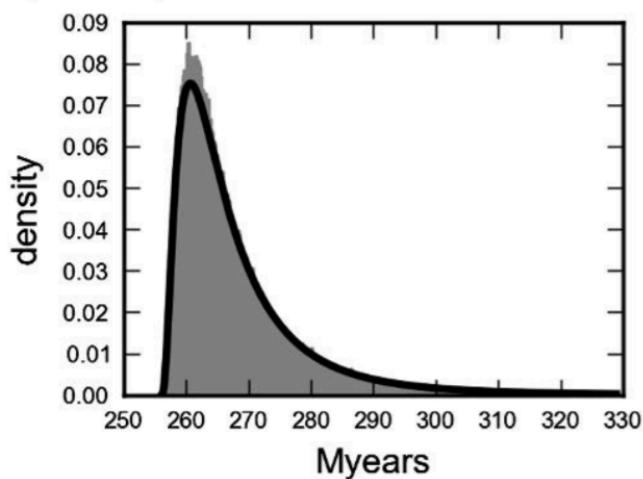


# Calibration nodes

---

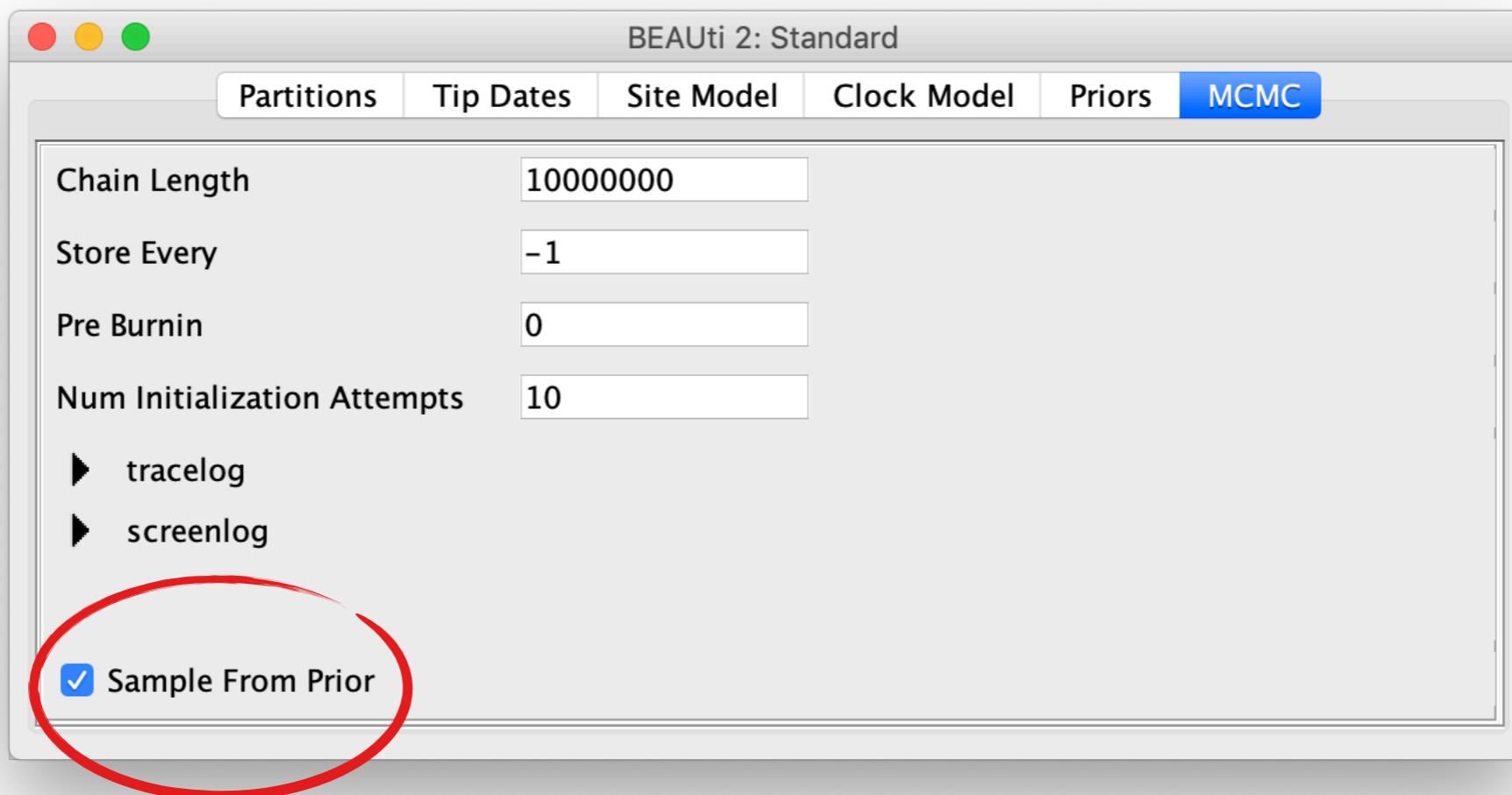
More calibration nodes mean more information to calibrate the clock

but is more always better?



Induced prior may be different to the prior you set!

# Setting priors best practice



- Sample or simulate from prior
- For key parameters plot prior and posterior together
- Mess with priors and see how sensitive/robust posterior is to prior perturbations

# On the program for today

---

- (1) What is phylodynamics?
- (2) Bayesian inference recap
- (3) BEAST2 introduction

**Tutorial:** Molecular clock dating (part i)

- (4) Molecular clock models

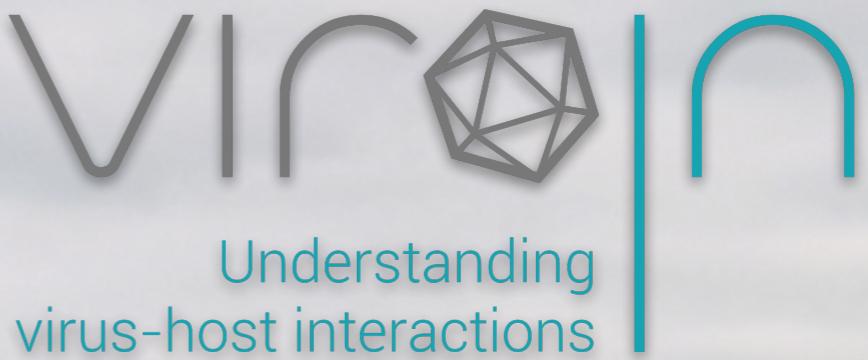
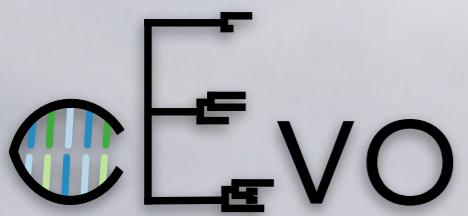
**Tutorial:** Molecular clock dating (part ii)

- (5) Setting priors

**Tutorial: Phylodynamics (part i)**

- (6) Tree priors

**Tutorial:** Phylodynamics (part ii)



0111010  
10100110  
10011010  
01111010  
01110011

Tree priors

Louis du Plessis

ETH zürich

DBSSE

# On the program for today

---

- (1) What is phylodynamics?
- (2) Bayesian inference recap
- (3) BEAST2 introduction

**Tutorial:** Molecular clock dating (part i)

- (4) Molecular clock models

**Tutorial:** Molecular clock dating (part ii)

- (5) Setting priors

**Tutorial:** Phylodynamics (part i)

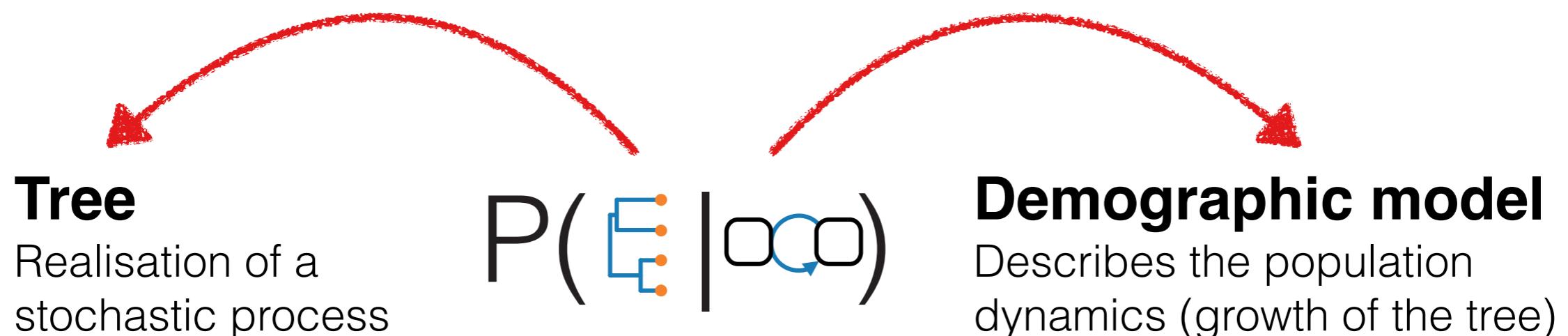
- (6) Tree priors**

**Tutorial:** Phylodynamics (part ii)



# Demographic model

- Describes the population/speciation dynamics
- How does the population demographics / species diversity change over time?

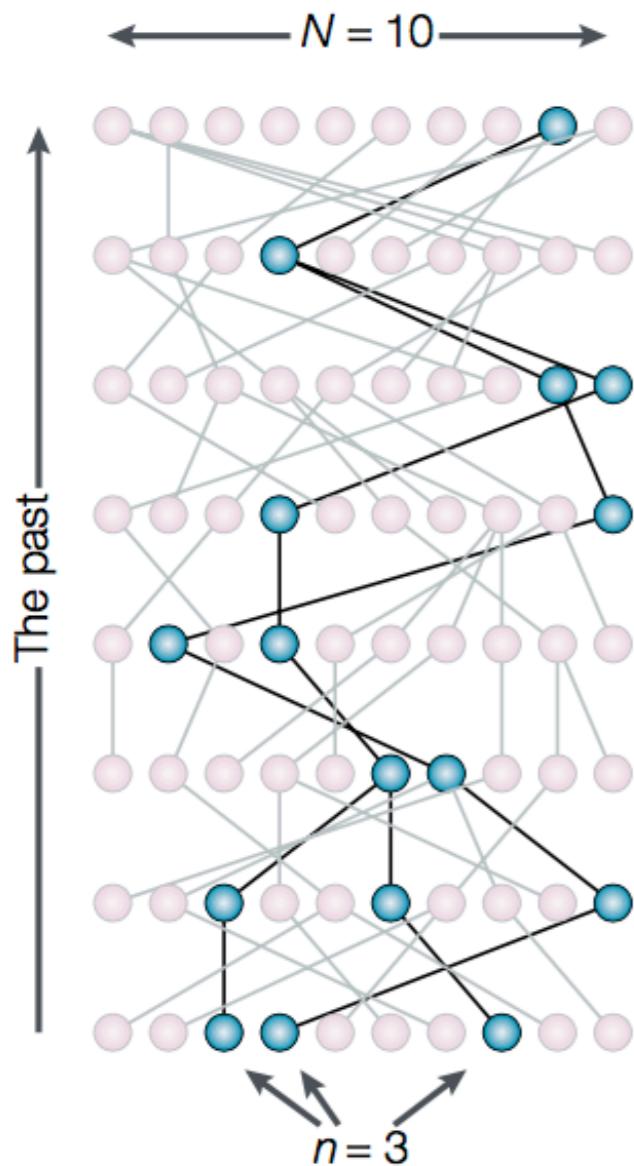


- How likely is the genealogy given a demographic model?
- Sometimes called a **tree prior**
- Usually a **coalescent** or **birth-death** model

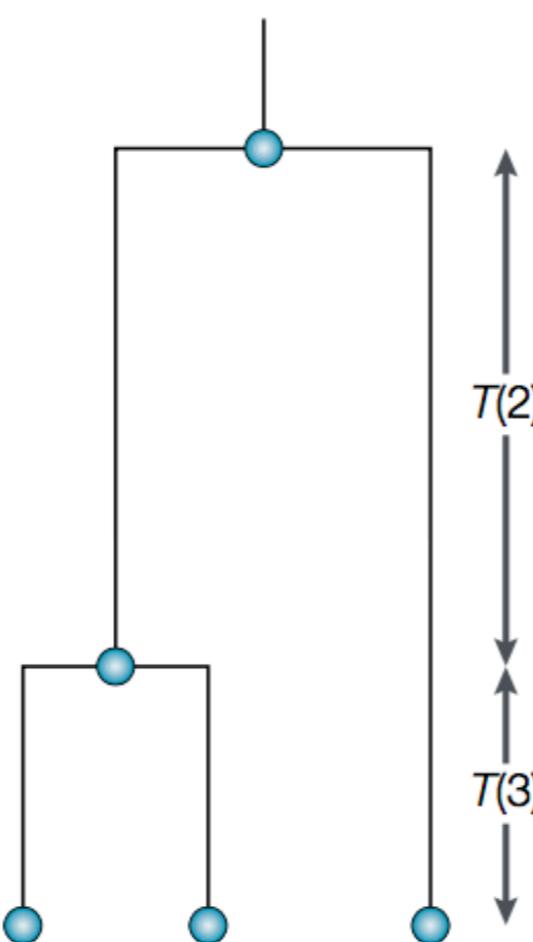


# Coalescent models

## Full population

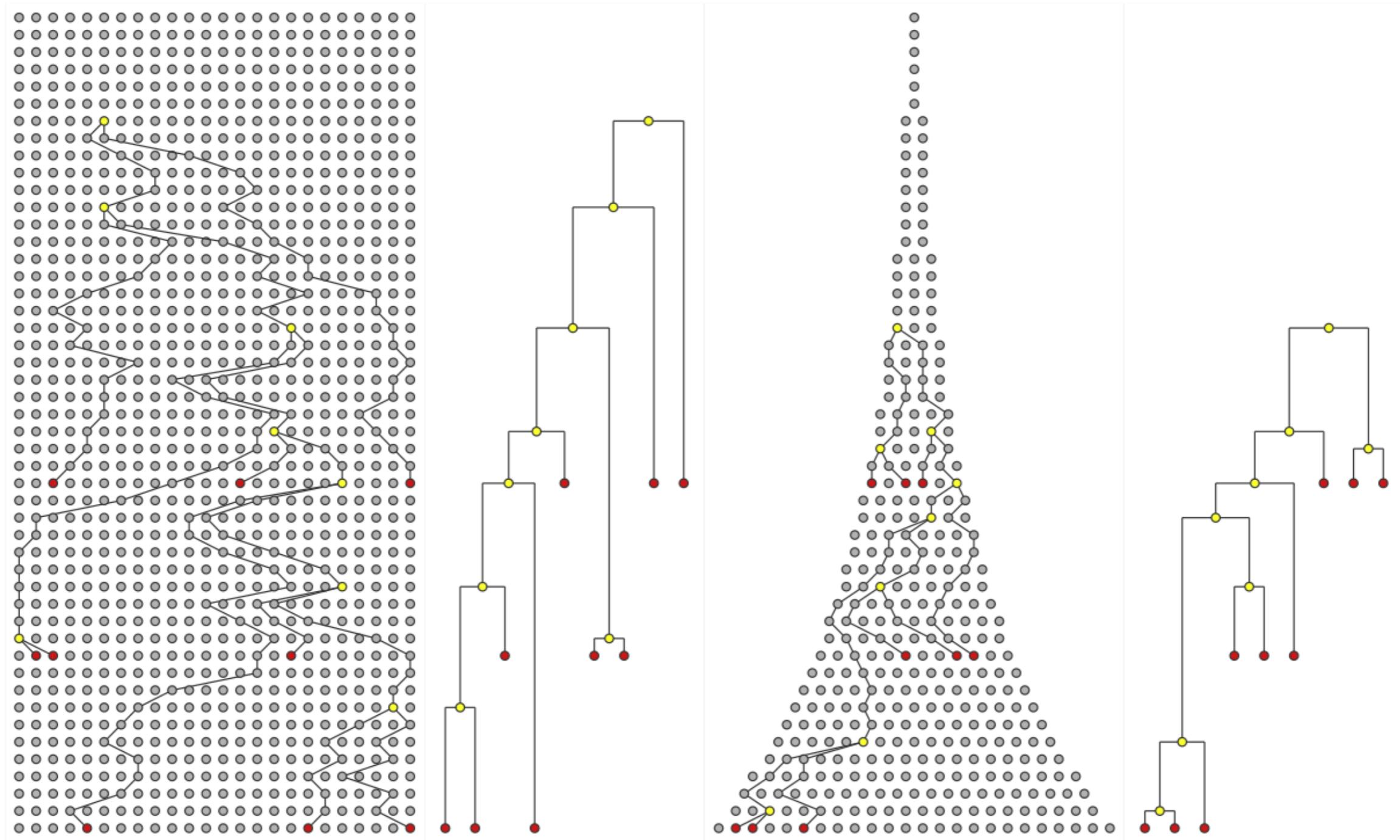


## Sampled tree



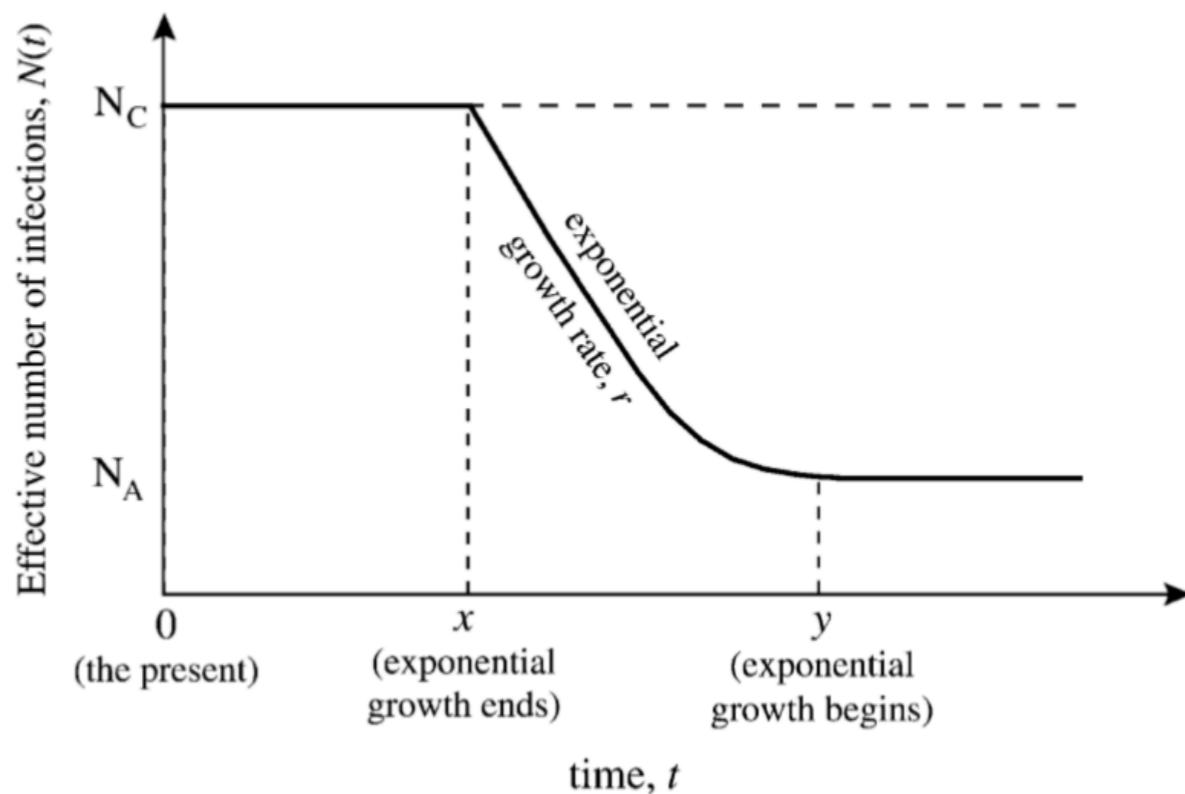
- Approximation to Wright-Fisher population dynamics (with large **N**)
- Discrete, non-overlapping generations, complete mixing
- Trace ancestry of **n** samples in a population of size **N**
- Given **N** it is easy to calculate the probability for **2** nodes to coalesce in time **t**
- Extend to the probability of **2** of **n** nodes coalescing in time **t**
- Calculate the probability of observing a given **tree** for a particular **N**  
→ estimate **N** (**N<sub>e</sub>** in practice)

# Changes in $N_e$ over time



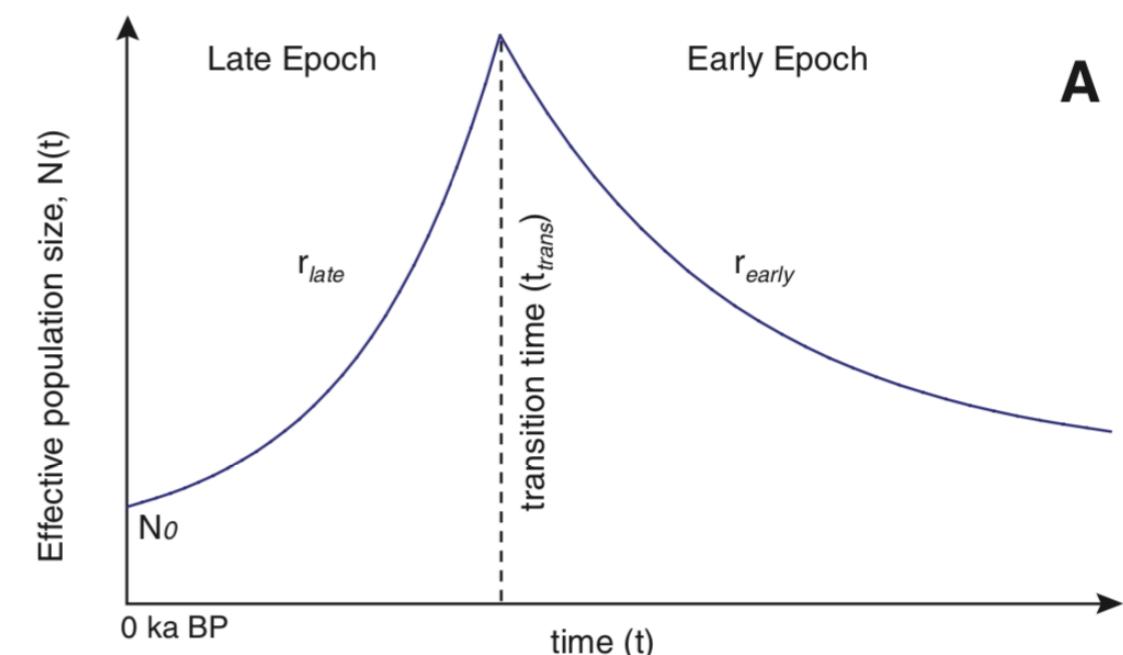
# Parametric models for $N_e(t)$

## Egyptian HCV



Constant-Exponential-Constant  
dynamics

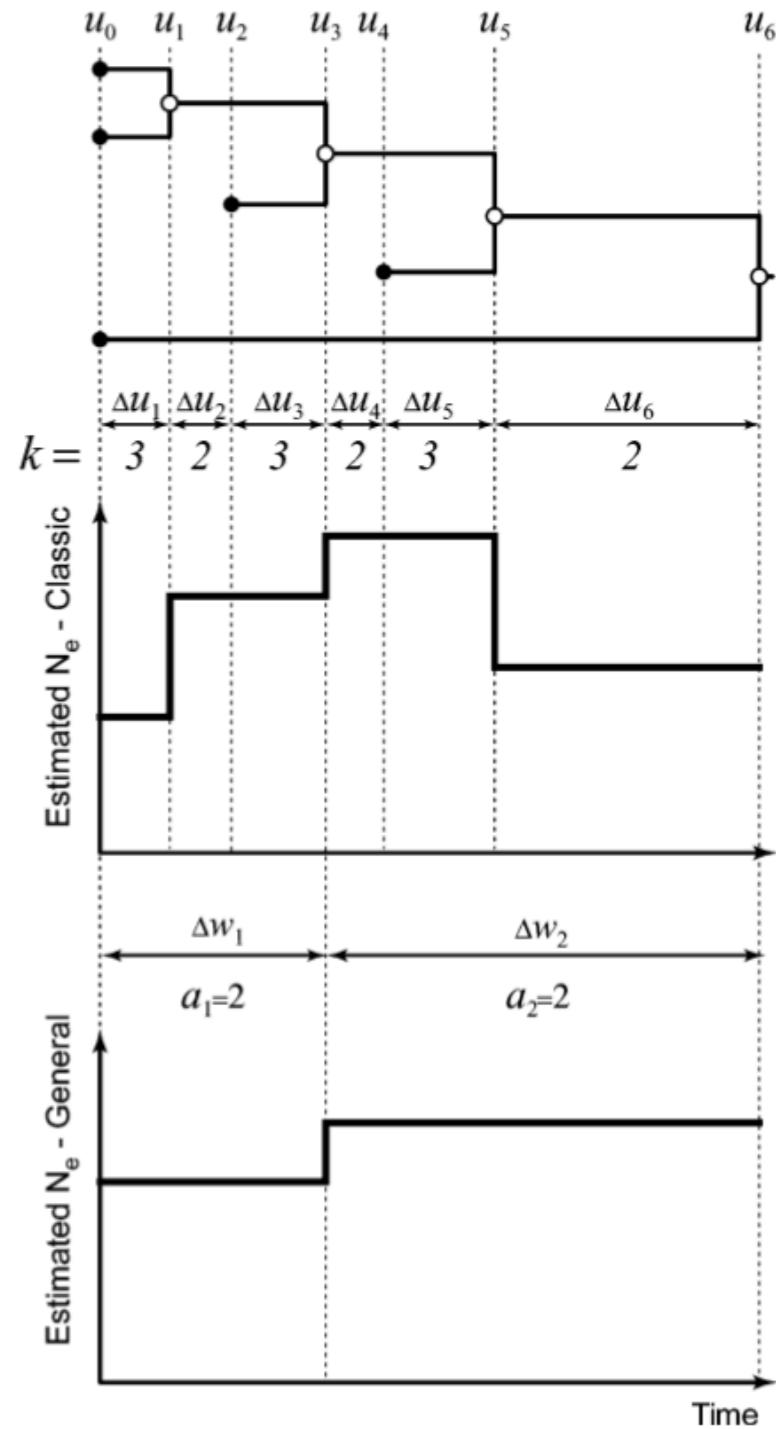
## Steppe Bison



Boom-Bust dynamics



# Skyline plots



## Classic skyline plot

- Piecewise constant  $N_e$
- Change-points at coalescence events
- Noisy estimate

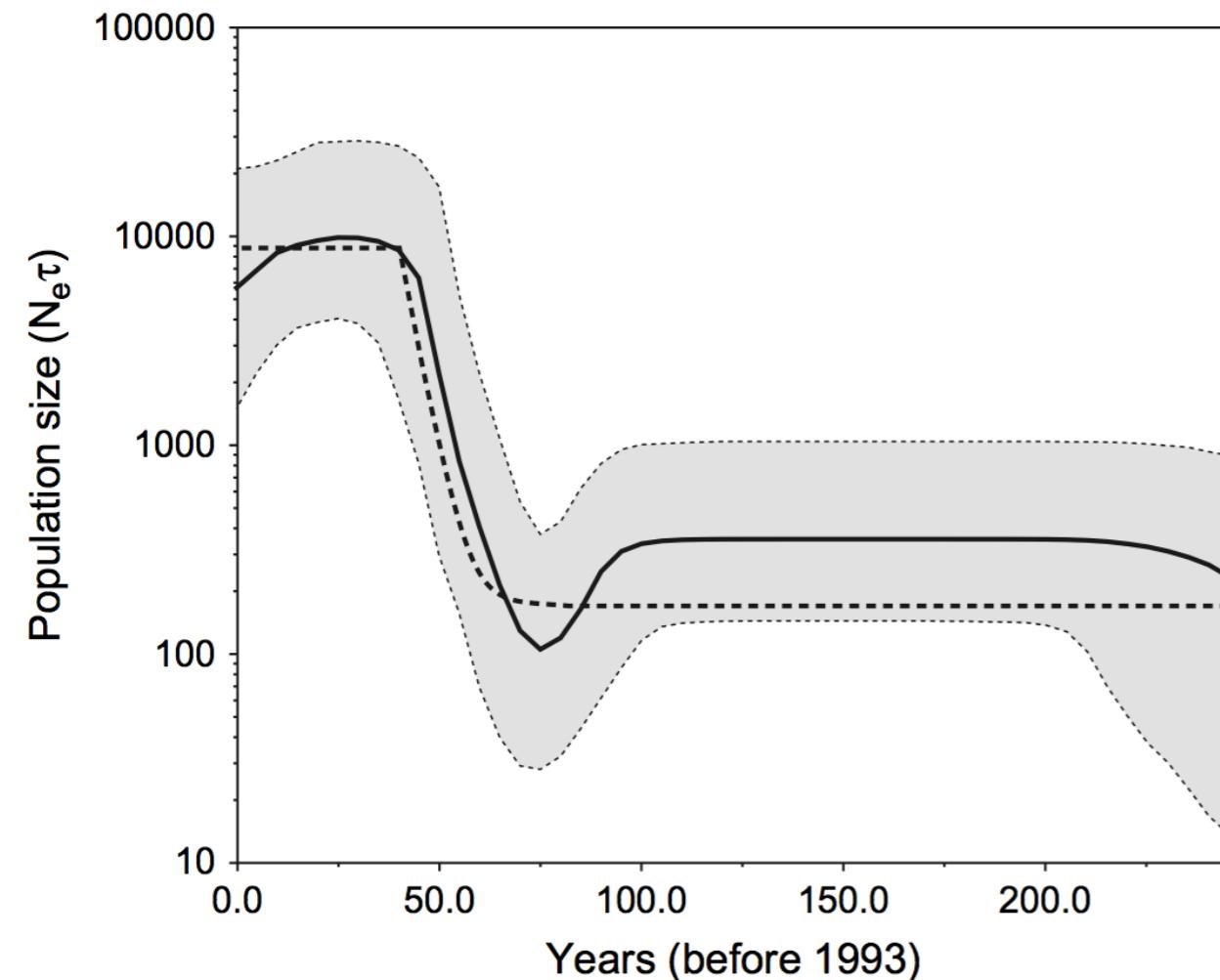
## Generalised skyline plot

- Group neighbouring segments
- Smoother estimate

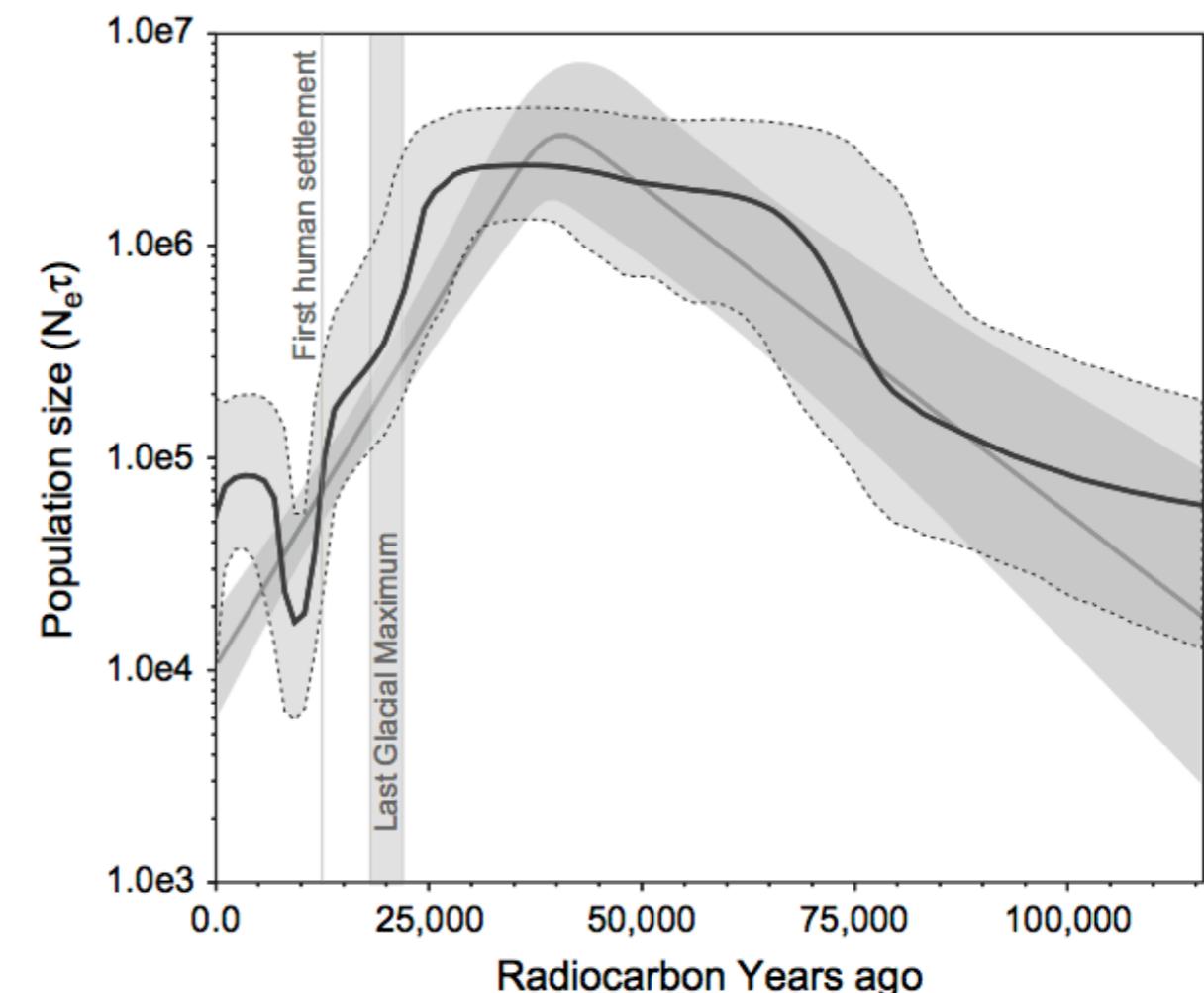


# Bayesian Skyline plot

## Egyptian HCV



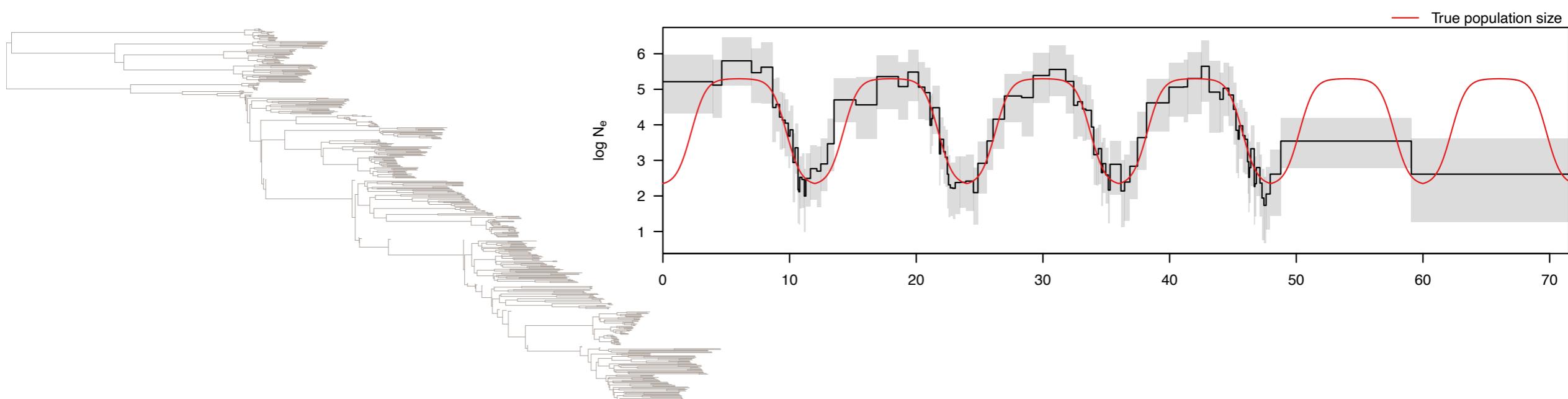
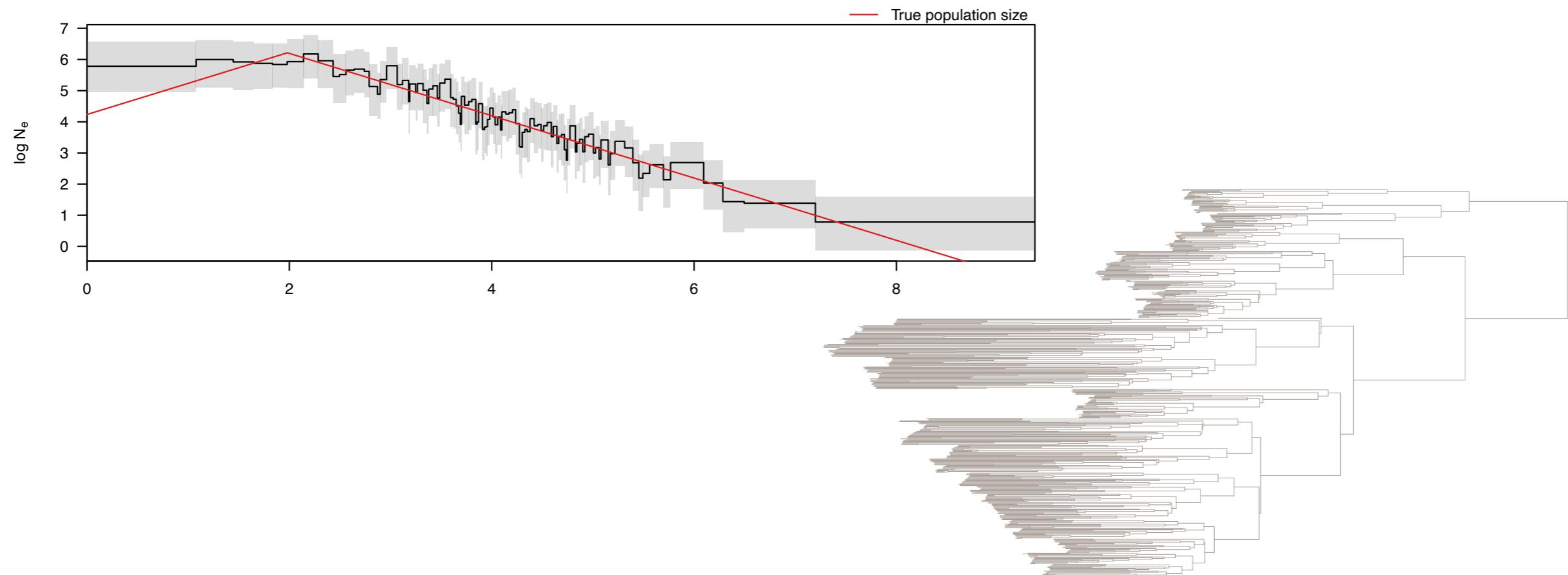
## Steppe Bison



- Recovers the same dynamics as the parametric model
- More flexible but also more uncertain

# How well does it work?

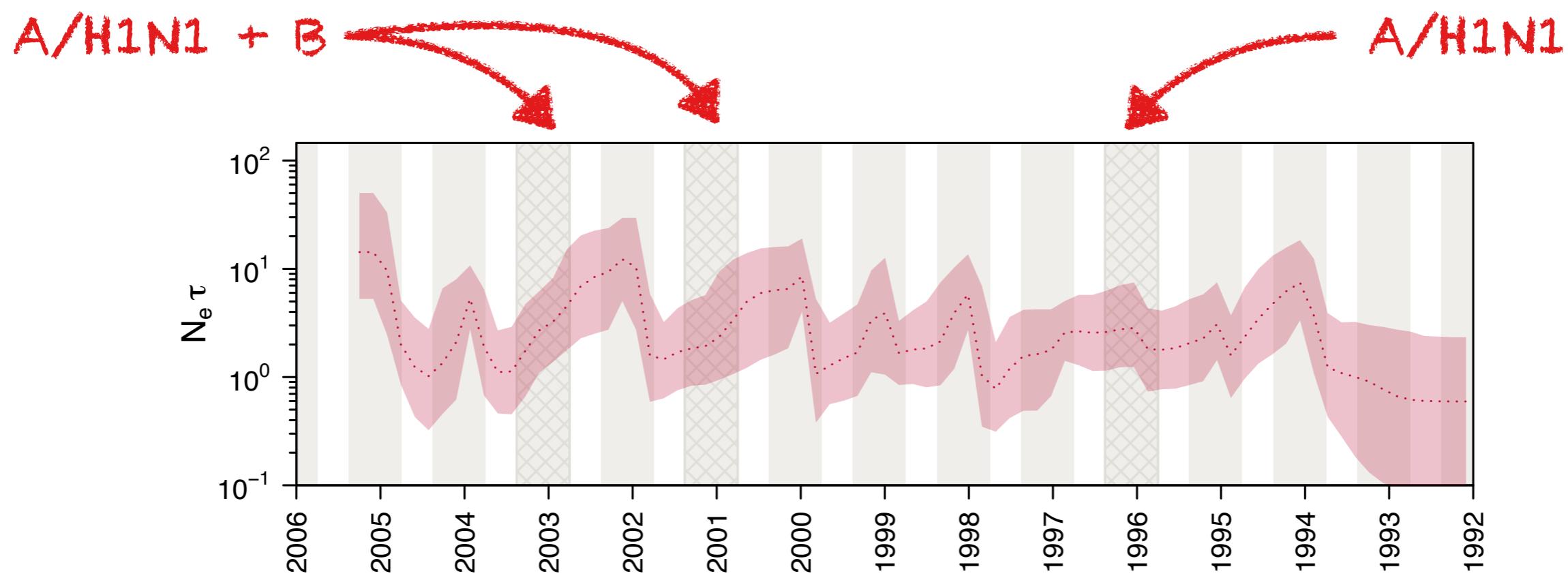
---





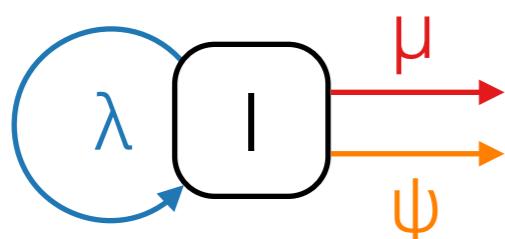
# Bayesian Skyline plot

## Influenza A H3N2 in New York



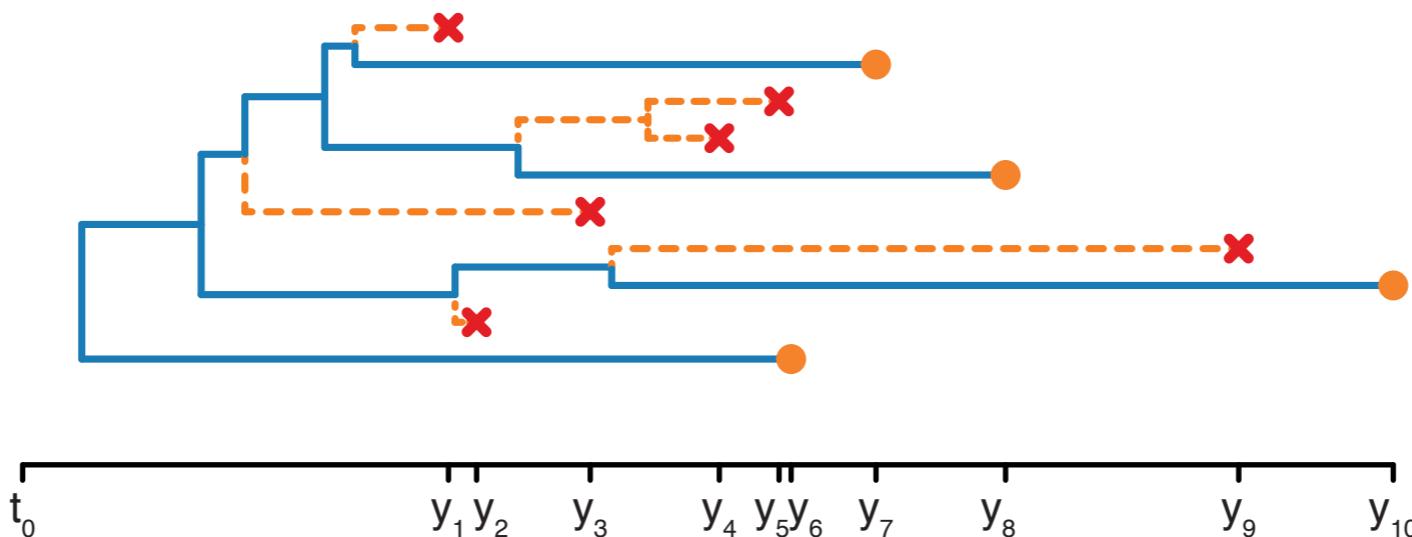


# Birth-death-sampling models



- $\lambda$  — birth/transmission rate (lineages added to **full** tree)
- $\mu$  — death/recovery rate (lineages removed from the **full** tree)
- $\psi$  — sampling rate (samples added to **sampled** tree)

- Forward-in-time branching process
- Events happen at different rates
  - infection/recovery
  - speciation/extinction
  - sampling/fossilization
  - ...



- Calculate the probability of a series of events at specific times to generate a (sampled) tree



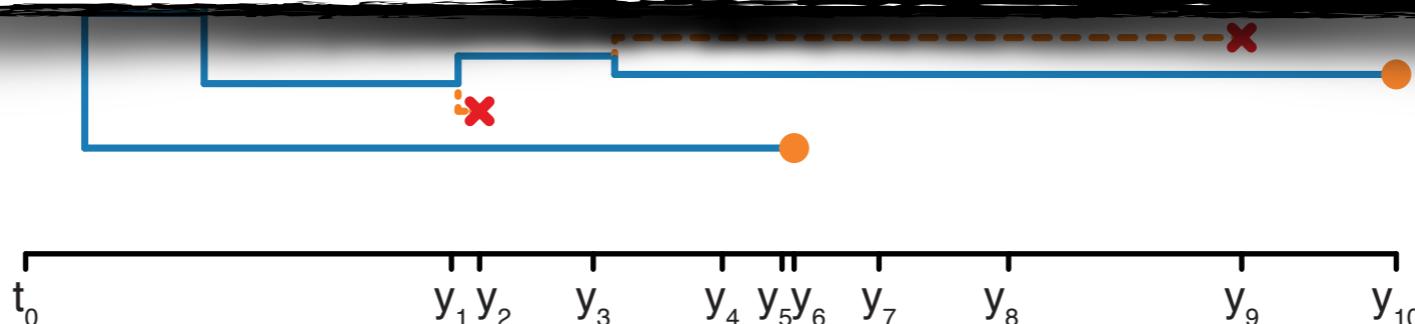
# Birth-death-sampling models

( $\lambda$ )

## Infectious disease parameterisation

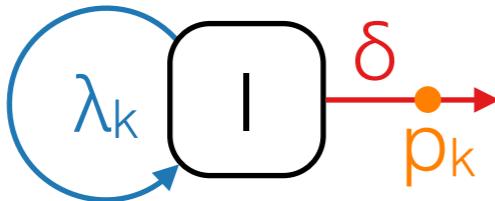
- $R_e = \lambda/(\mu + \psi)$  — Effective reproductive number  
(is it spreading or not?)
- $\delta = \mu + \psi$  — Becoming uninfected rate  
( $1/\delta$  = infectious period)
- $p = \psi/(\mu + \psi)$  — Sampling proportion  
(Proportion of sampled removals)

**More natural priors!**

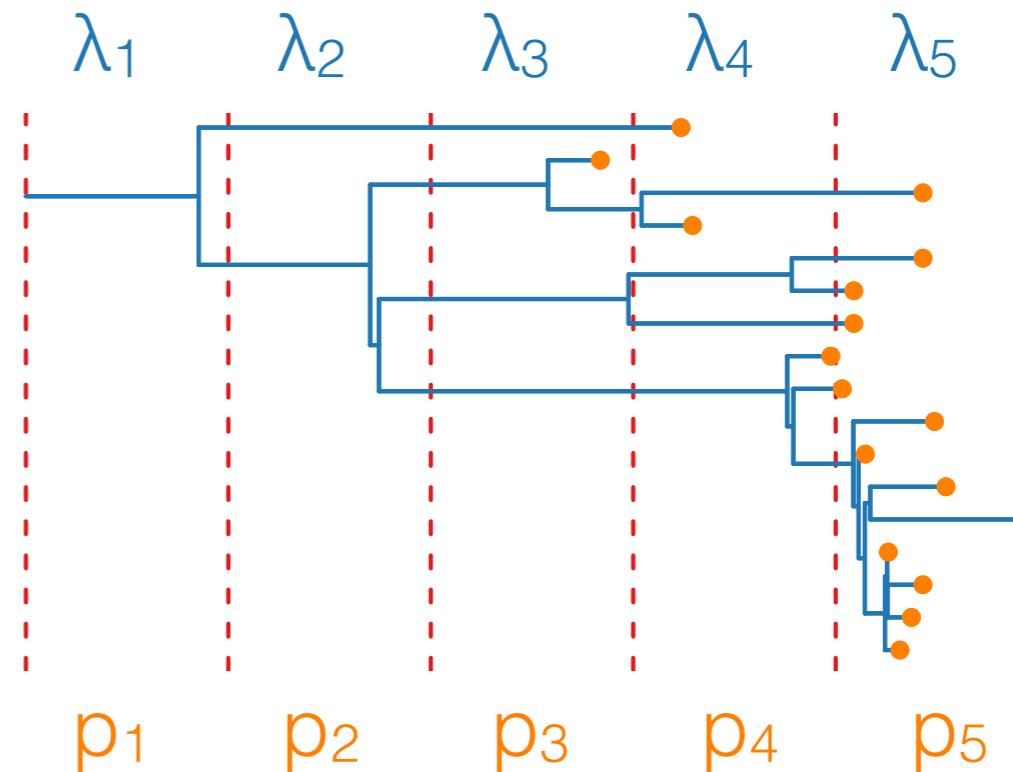


- Calculate the probability of a series of events at specific times to generate a (sampled) tree

# Birth-death skyline

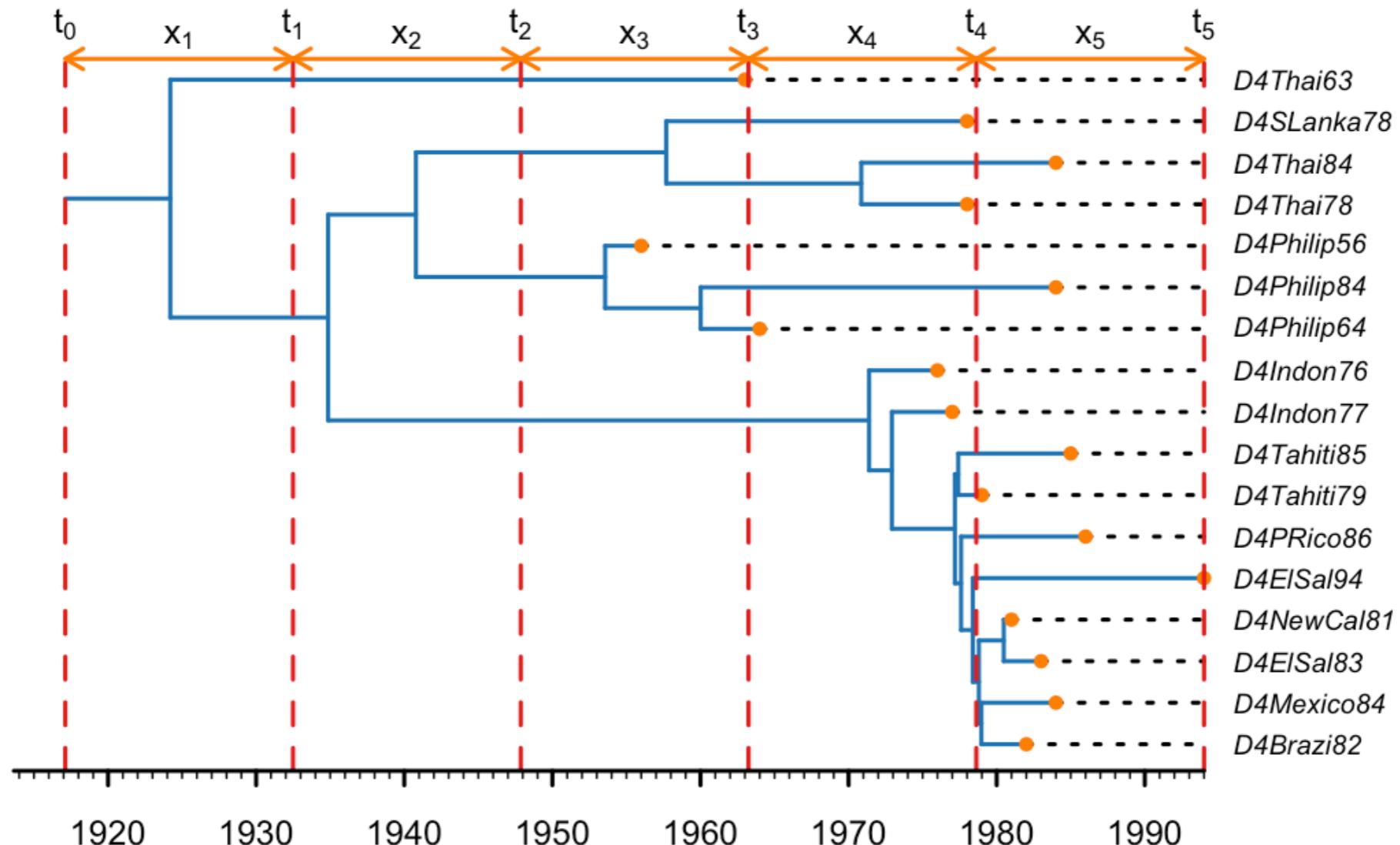


- $\lambda$  — infection rate
- $\delta$  — becoming-noninfectious rate
- $p$  — sampling proportion
- $t_0$  — Origin of the process



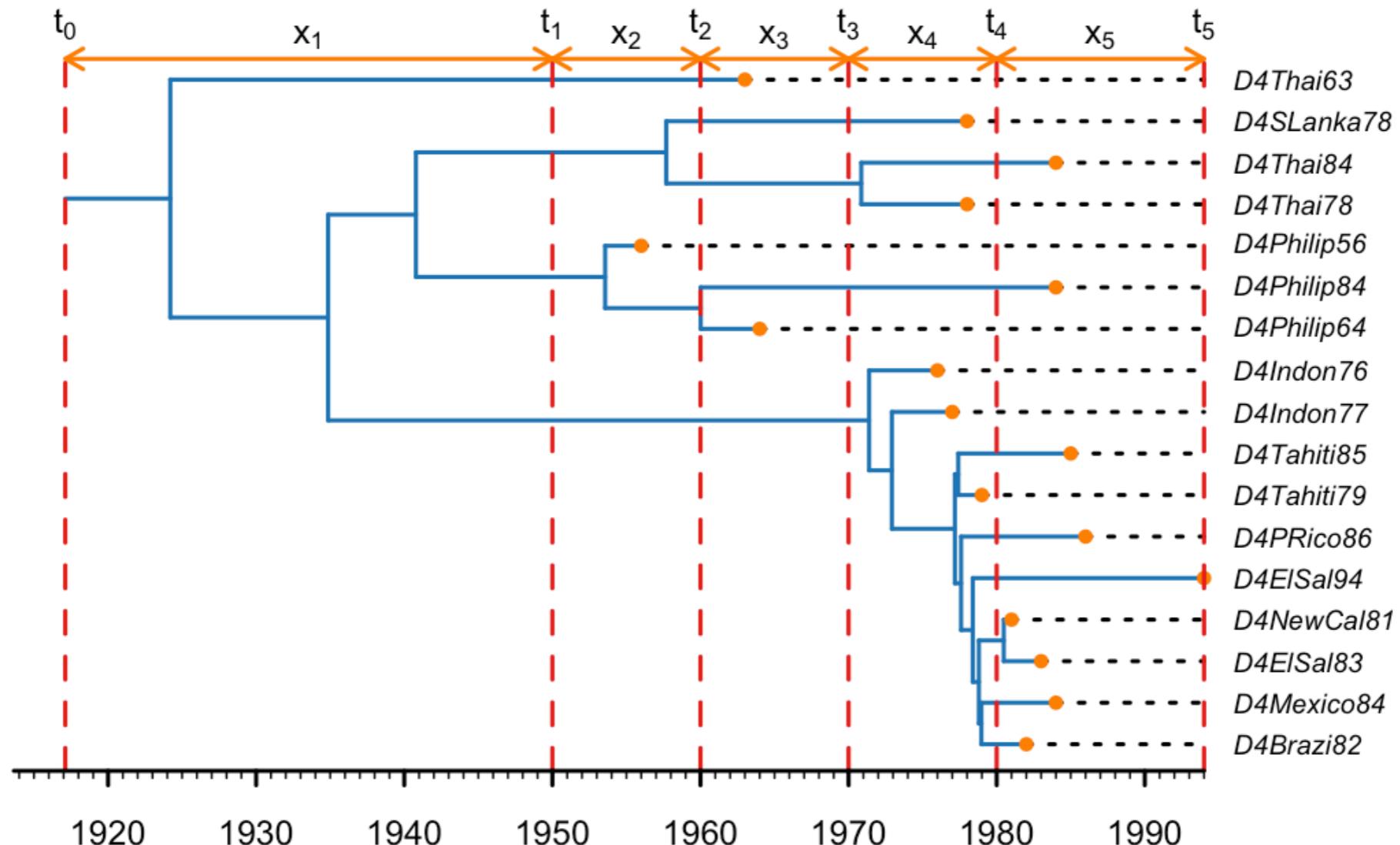
- Model parameters are rates that describe the growth of the tree
- Effective reproduction number =  $\lambda/\delta$   
(is it spreading or not?)
- Allow parameters to change through time
- Shifts in rates can be anywhere

# Birth-Death Skyline (BDSKY)



- Time-changing parameter can be **any** or **all** of the model rates (birth, death, sampling)

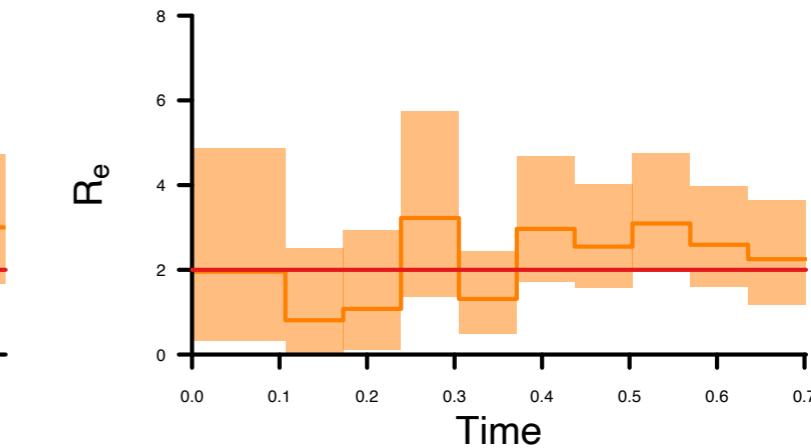
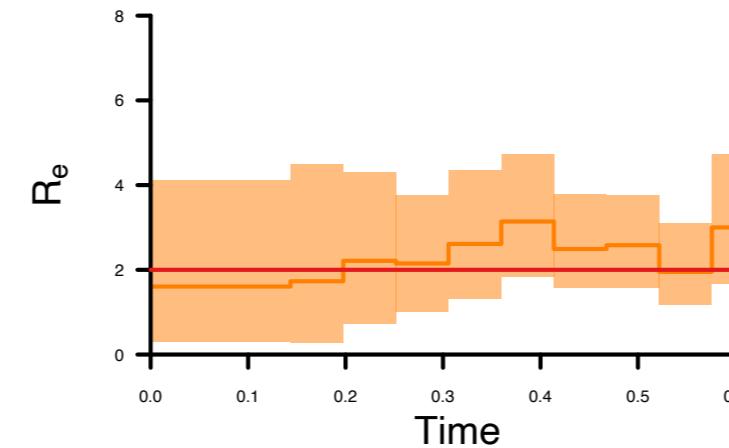
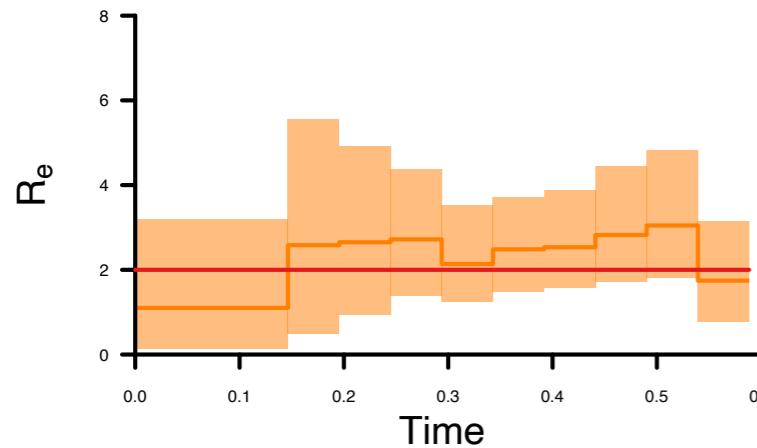
# Birth-Death Skyline (BDSKY)



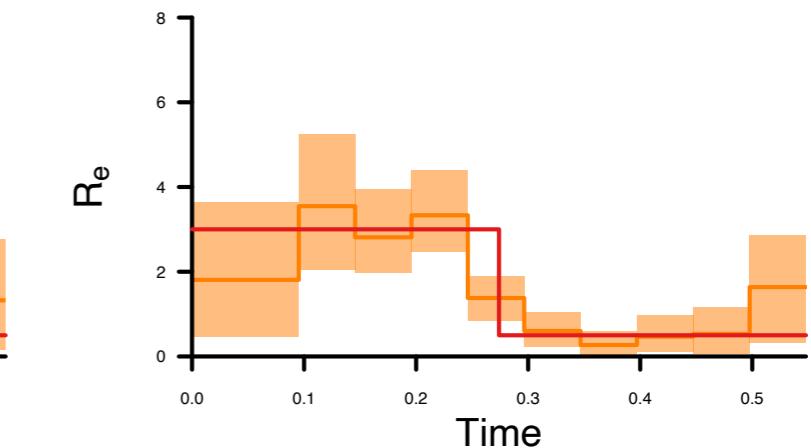
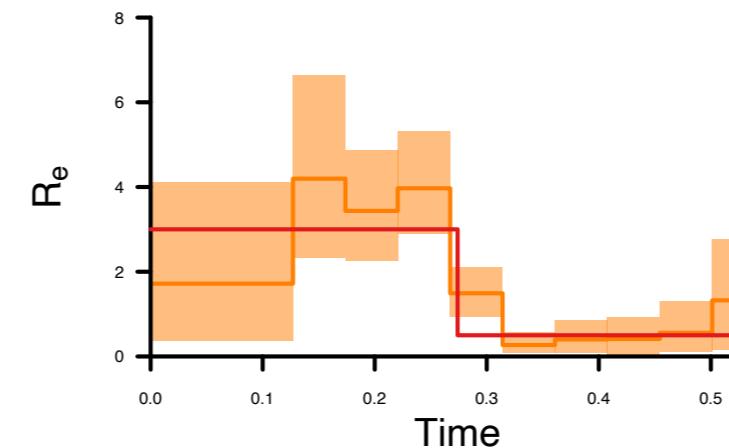
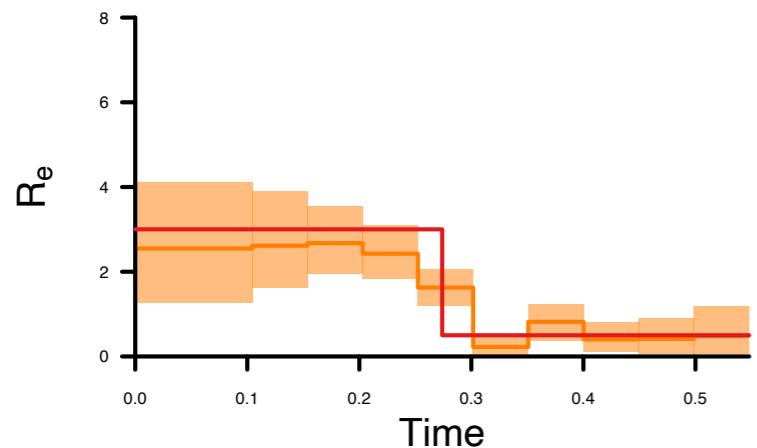
- Change-points can be anywhere between origin and present
- More difficult to set up XML if not equally-spaced

# How well does it work?

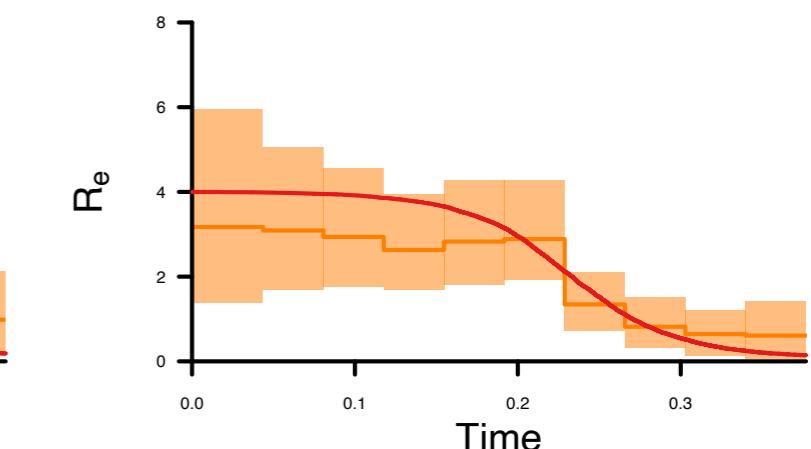
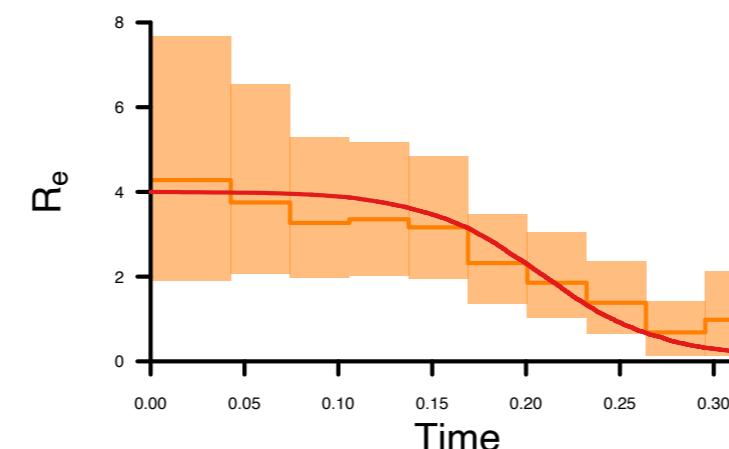
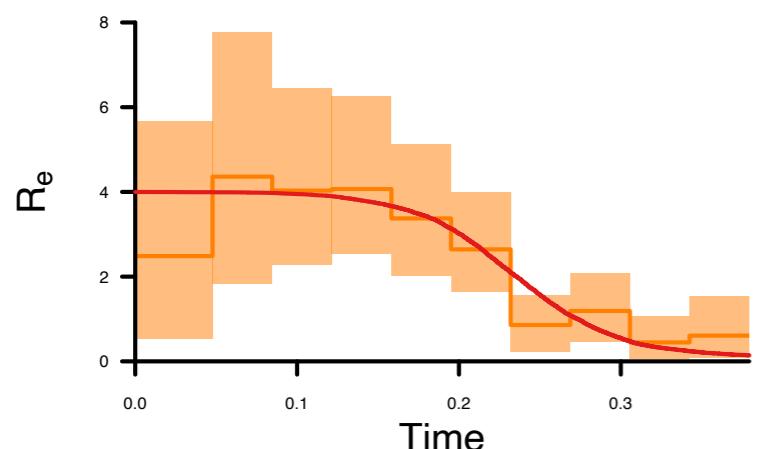
Constant



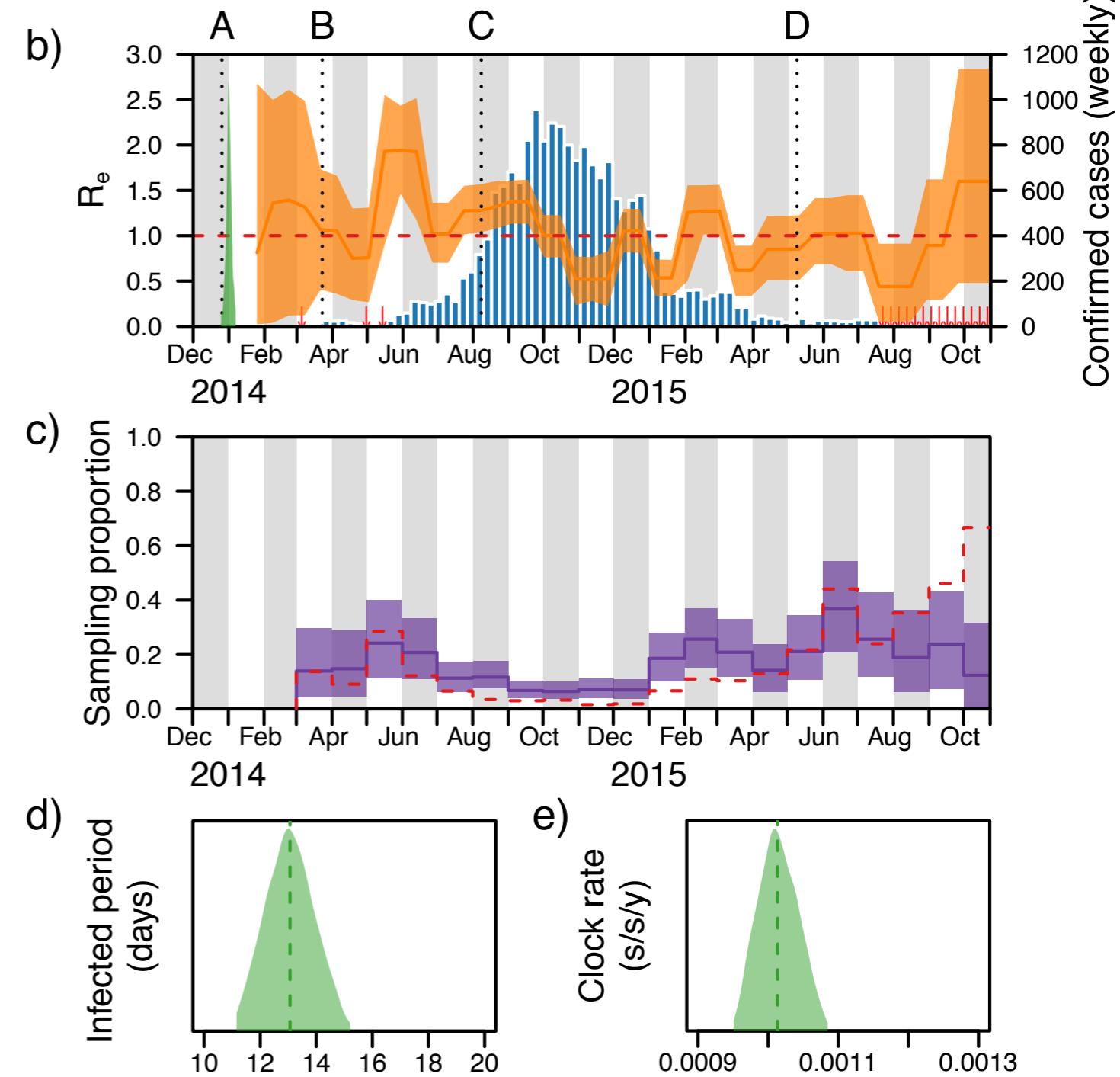
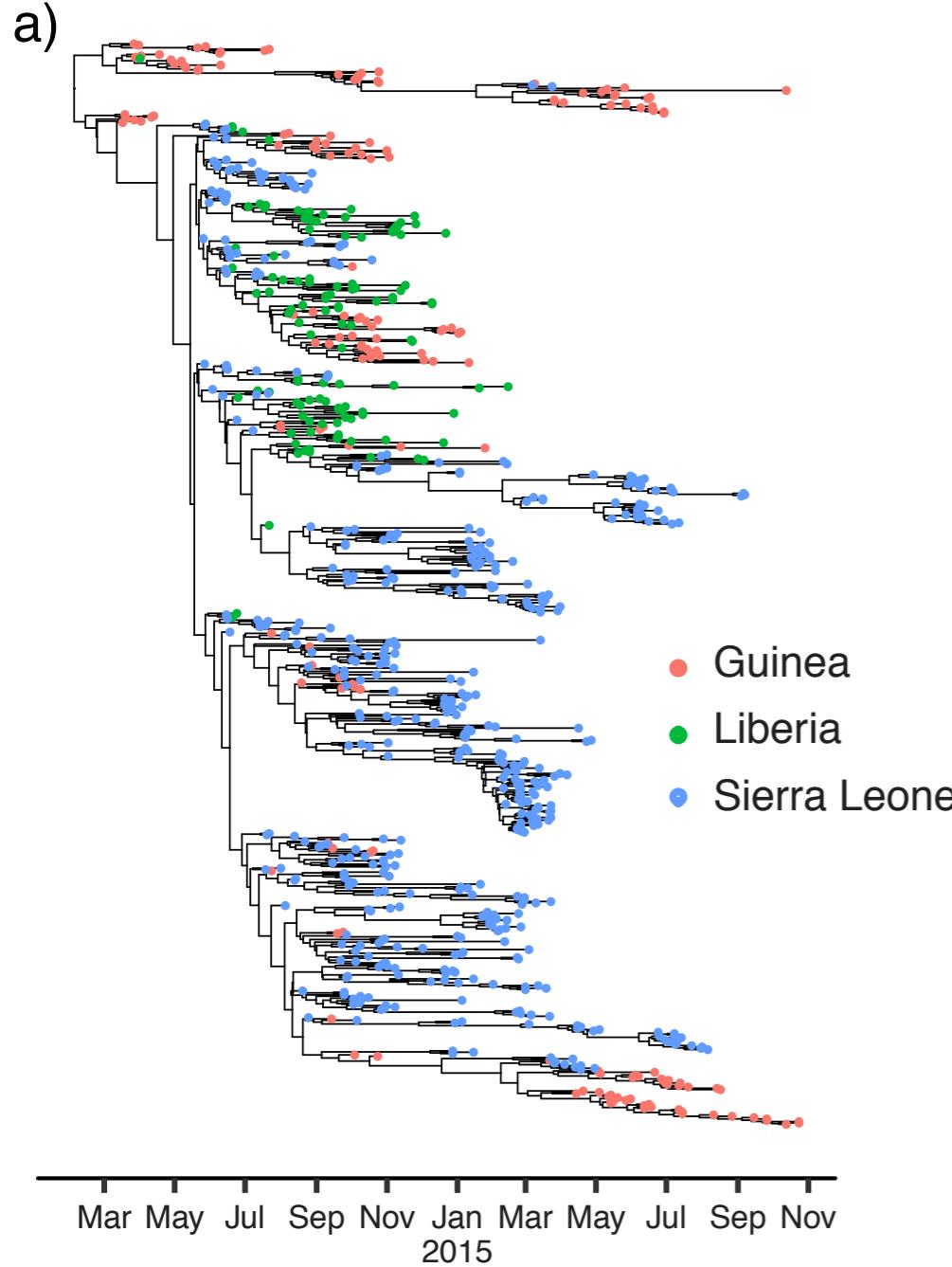
Piecewise constant



SEIR model



# EBOV in West Africa





# Demographic model



## Tree

Realisation of a stochastic process

$$P(F | D)$$

## Demographic model

Describes the population dynamics (growth of the tree)

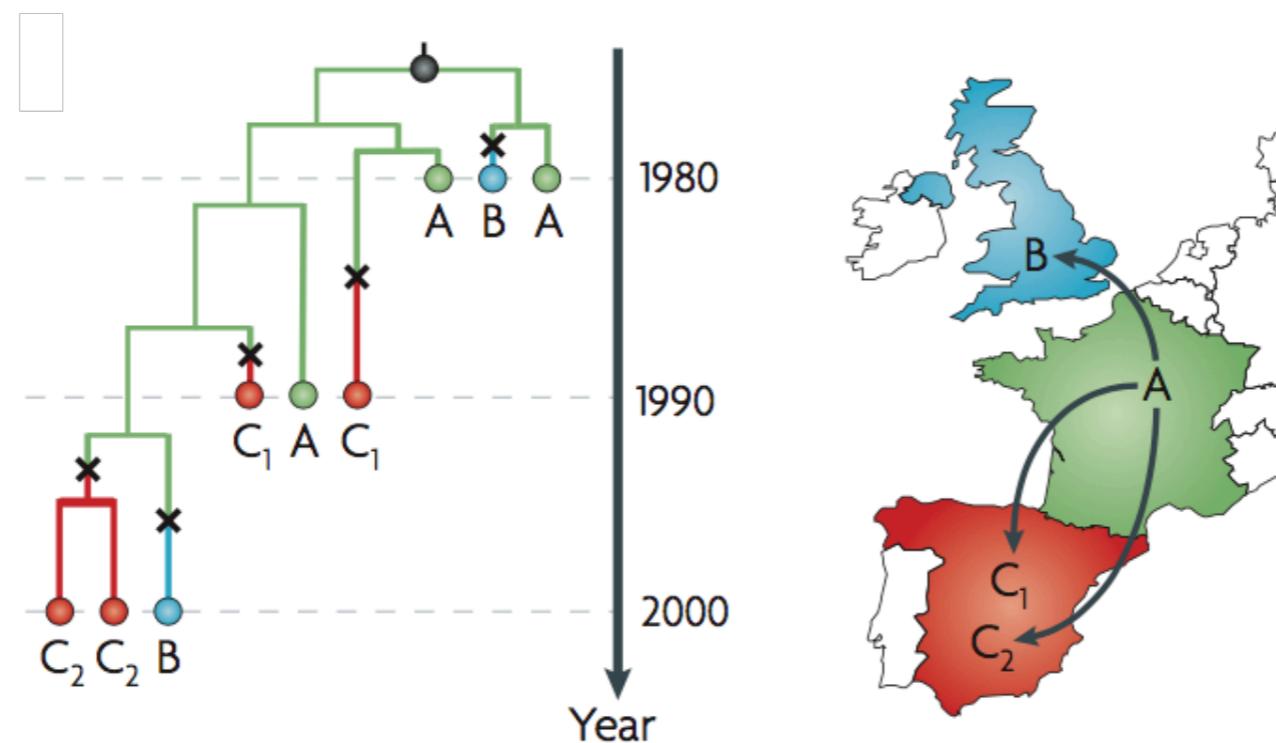
- **Coalescent:**

Given  $n$  sampling times and an estimate for  $N_e(t)$ , out of all the ways we can connect the samples, what is the probability of the current tree?

- **Birth-death:**

Given an estimate for the **origin** time, **birth**, **death** and **sampling** rates, if we simulate a tree forward-in-time from the origin to the time of the most recent sample, out of all the trees with  $n$  samples, what is the probability of the current tree?

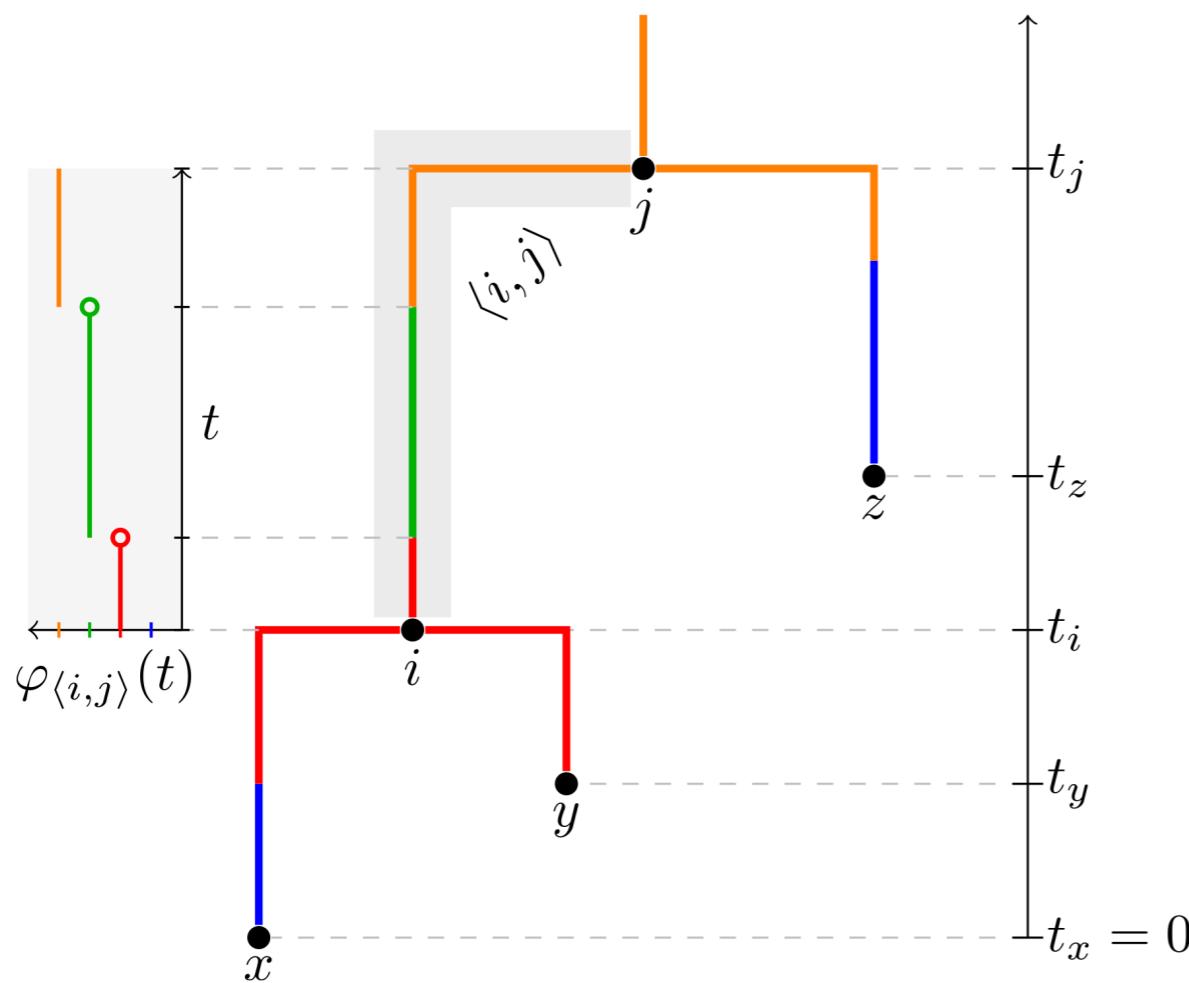
# Phylogeography



- "Migration" model assumes a migration model independent of the tree prior (treats migrations like substitutions - extra alignment with one site for location)
- True structured models (structured coalescent, multi type birth-death model) model it as part of the tree-prior

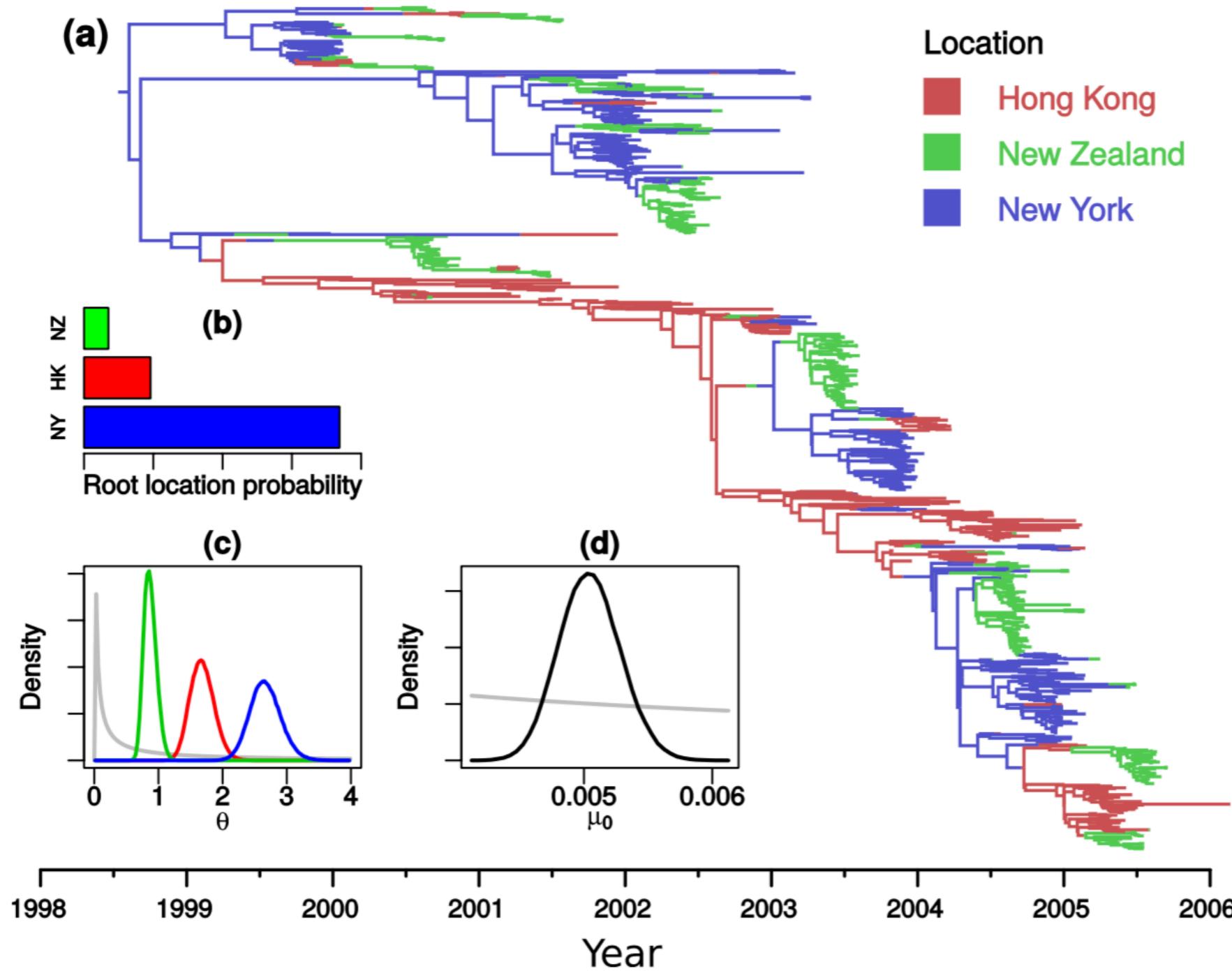


# Structured models



- Explicitly model different types of subpopulations or locations
- Different types may be associated with different epidemiological dynamics
- Can also model migration between types

# Seasonal Influenza H3N2 phylogeography



# Thank you for listening!

Slides are my own, but some slides were inspired by (or copied from) slides by **Alexei Drummond**, **David Rasmussen**, **Tanja Stadler**, **Simon Ho** and **Oliver Pybus**

