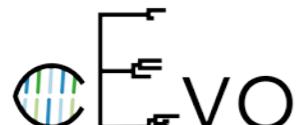


Taming the BEAST

Bayesian evolutionary analysis by sampling trees



**ACACACCTACAGACTTACAGACCC
TCACACCTACACACCCACAGACTT
TCAGACTTTCACACCTTCAGACCT
ACAGACTTTCAGACTTTCAGACCC
TCACACCTACACACCCACAGACTT
TCAGACTTTCACACCTTCAGACCT
TCACACCTACACACCCACAGACTT
TCAGACTTTCACACCTTCAGACCT**



Computational Evolution

Louis du Plessis and Carsten Magnus



D-BSSE
Department of Biosystems
Science and Engineering

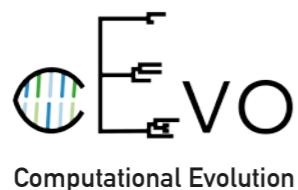
Taming the BEAST

Bayesian evolutionary analysis by sampling trees



1. The problem

ACACACCTACAGACTTACAGACCC
TCACACCTACACACCCCACAGACTT
TCAGACTTTCACACCCCTTCAAGACCT
ACAGACTTTCAAGACTTTCAAGACCC
TCACACCTACACACCCCACAGACTT
TCAGACTTTCACACCCCTTCAAGACCT
TCACACCTACACACCCCACAGACTT
TCAGACTTTCACACCCCTTCAAGACCT



Computational Evolution

Louis du Plessis and Carsten Magnus



D-BSSE
Department of Biosystems
Science and Engineering

Taming the BEAST

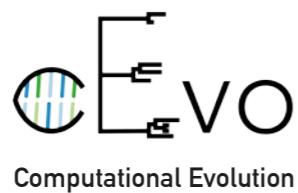
Bayesian evolutionary analysis by sampling trees



ACACACCTACAGACTTACAGACCC
TCACACCTACACACCCCACAGACTT
TCACACCTACACACCCCACAGACTT
ACAGACTTTCAGACTTTCAGACCC
TCACACCTACACACCCCACAGACTT
TCAGACTTTCACACCTTCAGACCT
TCACACCTACACACCCCACAGACTT
TCAGACTTTCACACCTTCAGACCT

1. The problem

2. What goes into a BEAST model?



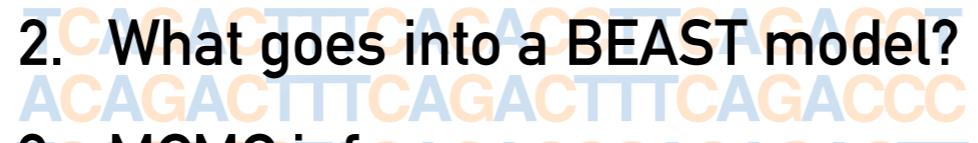
Louis du Plessis and Carsten Magnus

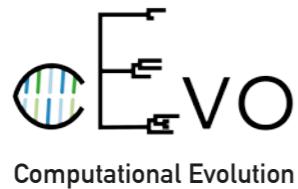


Taming the BEAST

Bayesian evolutionary analysis by sampling trees



1. The problem
A sequence alignment diagram showing multiple DNA or RNA sequences. The top sequence is ACACACCTACAGACTTACAGACCC. Below it are several other sequences of varying lengths, primarily in shades of blue and orange, representing different parts of the same molecule or related molecules.
2. What goes into a BEAST model?
A sequence alignment diagram showing multiple DNA or RNA sequences. The top sequence is ACAGACTTT CAGACTTT CAGACCC. Below it are several other sequences of varying lengths, primarily in shades of blue and orange, representing different parts of the same molecule or related molecules.
3. MCMC inference
A sequence alignment diagram showing multiple DNA or RNA sequences. The top sequence is CCCCACAGACTT. Below it are several other sequences of varying lengths, primarily in shades of blue and orange, representing different parts of the same molecule or related molecules.



Computational Evolution

Louis du Plessis and Carsten Magnus



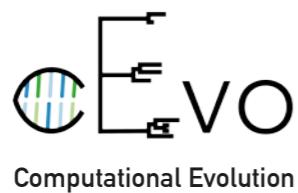
D-BSSE
Department of Biosystems
Science and Engineering

Taming the BEAST

Bayesian evolutionary analysis by sampling trees



1. The problem
2. What goes into a BEAST model?
3. MCMC inference
4. Bayesian evolutionary analysis
by sampling trees



Computational Evolution

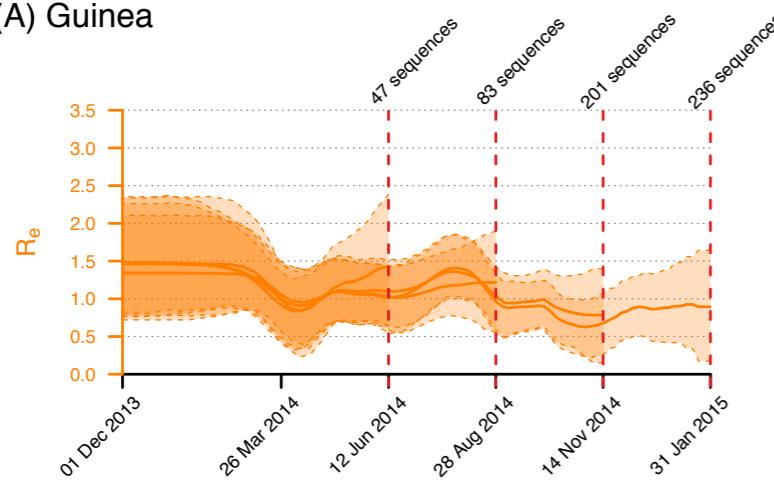
Louis du Plessis and Carsten Magnus



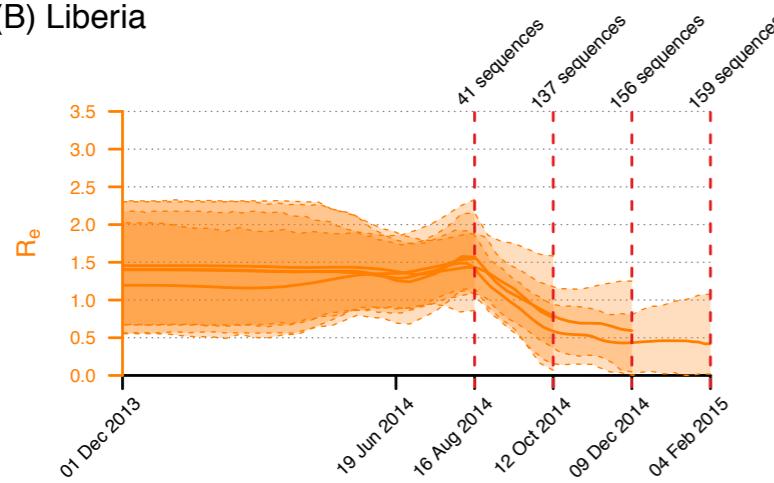
D-BSSE
Department of Biosystems
Science and Engineering

Myself

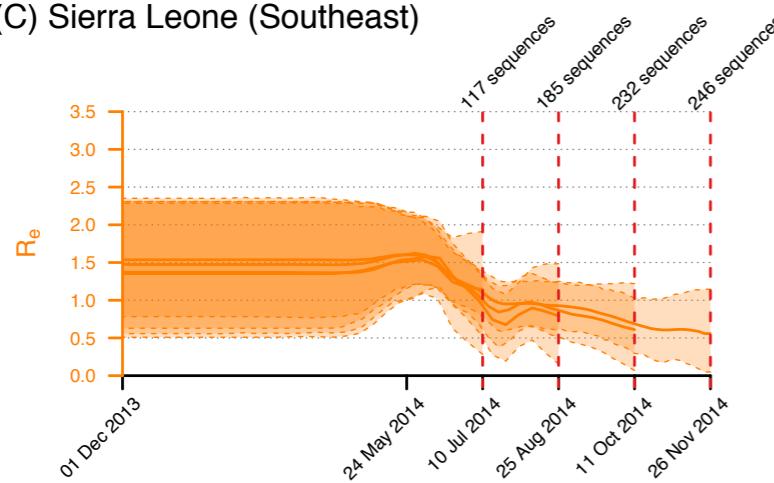
(A) Guinea



(B) Liberia



(C) Sierra Leone (Southeast)



Real-time phylodynamics

- Viral data sampled over the course of an outbreak

Model biases and limitations

- Simulated data

Everyone else

- Evolution of HIV-1, tuberculosis, leprosy, Influenza etc.
- Evolution of drug resistance
- Host-parasite co-evolution
- Disease spread by migration/travel
- Species delimitation
- Bird taxonomy
- Plant diversity and dispersal
- Adaptive radiations and diversification rates
- Time calibration using fossils
- Ancient DNA
- Influence of environmental factors on phylogenetic diversity
- Total evidence dating

...etc

We all have one thing in common...



All of us use genomic sequencing data
to answer questions in **BEAST**

We all have one thing in common...



All of us use genomic sequencing data
to answer questions in **BEAST**

Bayesian inference recap

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

Bayesian inference recap

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

The diagram illustrates the components of Bayes' theorem. A red curved arrow labeled "Likelihood" points from the term $P(\text{data} \mid \text{model})$ in the numerator to the term $P(\text{data})$ in the denominator. Another red curved arrow labeled "Prior" points from the term $P(\text{model})$ in the numerator to the same term in the denominator. A straight red arrow labeled "Posterior" points from the term $P(\text{model})$ in the numerator to the result of the division in the denominator.

Bayesian inference recap

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

The diagram illustrates the components of the Bayesian formula:

- Likelihood:** An arrow points from the term $P(\text{data} \mid \text{model})$ in the numerator to the word "Likelihood".
- Prior:** An arrow points from the term $P(\text{model})$ in the numerator to the word "Prior".
- Marginal Likelihood of the data:** An arrow points from the term $P(\text{data})$ in the denominator to the words "Marginal Likelihood of the data".
- Posterior:** An arrow points from the final result $P(\text{model} \mid \text{data})$ to the word "Posterior".

Bayesian inference recap

Data

- One or more alignments of sequencing data
- Sampled at one or many time points
- Timescale of days to millions of years
- Realisation of a stochastic process

Posterior

Prior

Marginal
Likelihood
of the data

Bayesian inference recap

Data

- One or more alignments of sequencing data
- Sampled at one or many time points
- Timescale of days to millions of years
- Realisation of a stochastic process

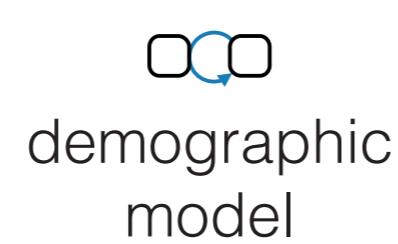
Model

- Description of the process that generated the data
- Parameters are random variables
- We may only be interested in some of the parameters
- Still need a prior for all the model parameters!

What goes into a BEAST model?



genealogy



demographic
model



site model



molecular clock
model

What goes into a BEAST model?



genealogy



demographic
model



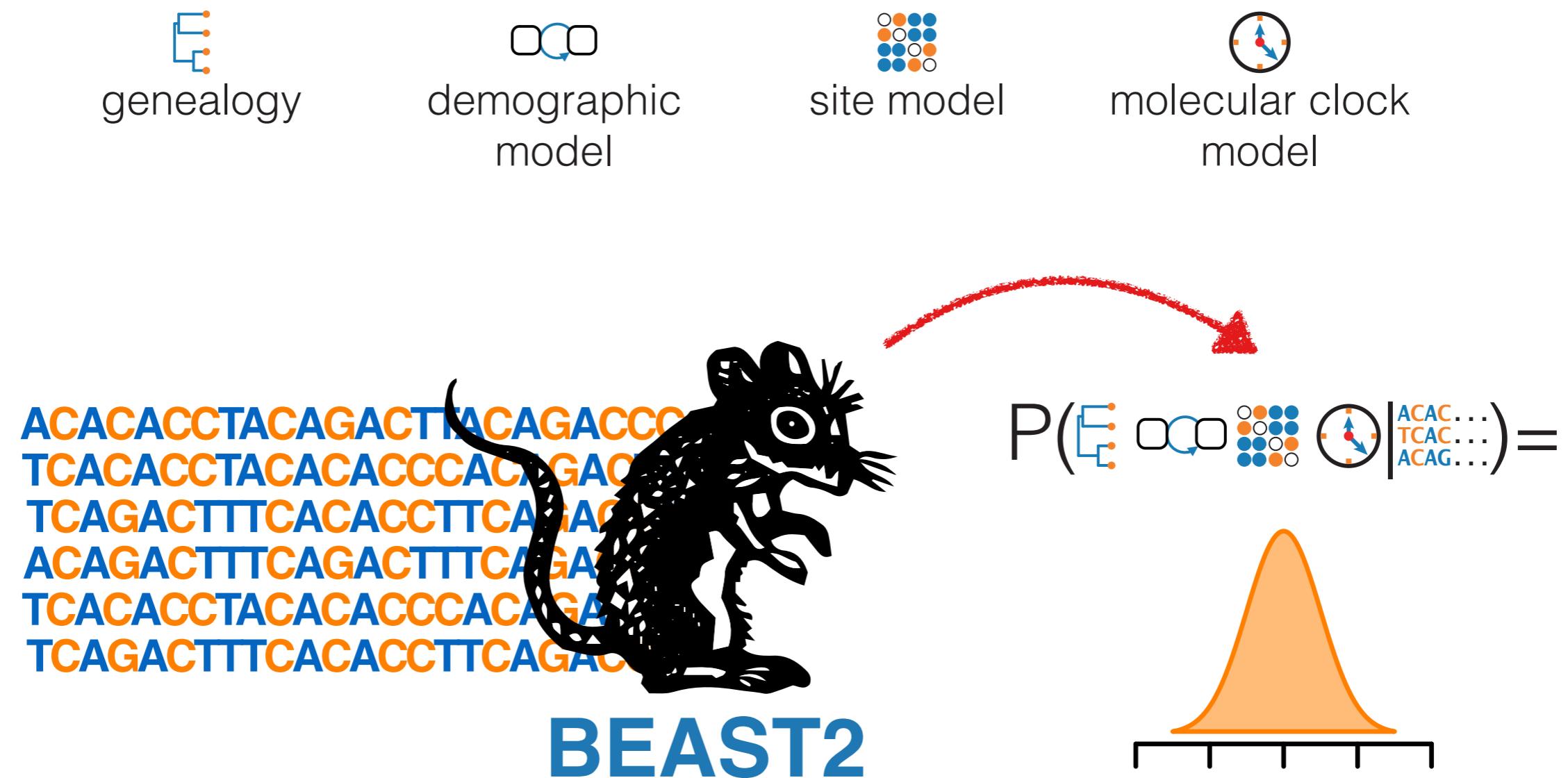
site model



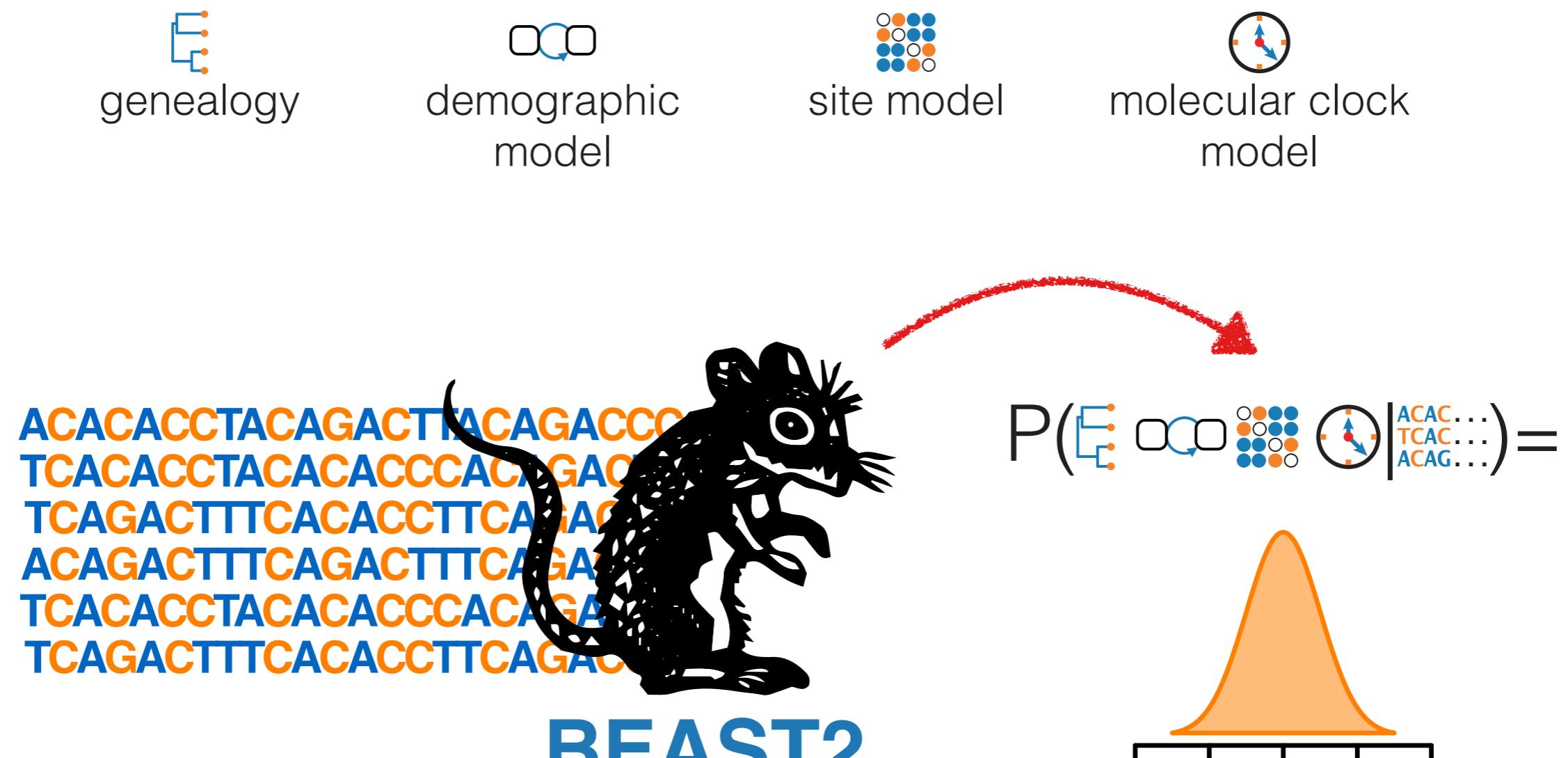
molecular clock
model



What goes into a BEAST model?



What goes into a BEAST model?



Estimate posterior distributions!



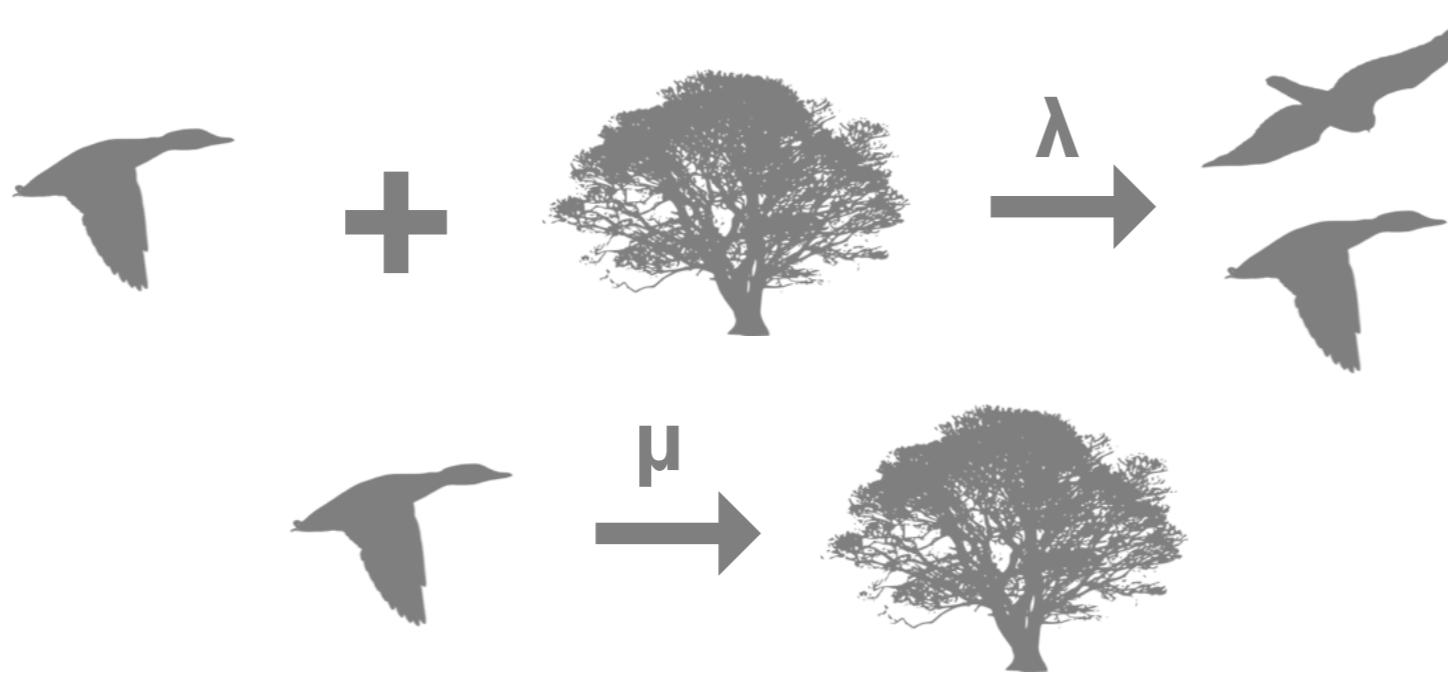
Demographic model

- Description of the population dynamics
 - How does the population grow over time?
 - Infectious diseases - transmission/recovery
 - Species - speciation/extinction
 - Migration events
 - Different classes of individuals
 - Host/vector dynamics
- ...etc**



Demographic model

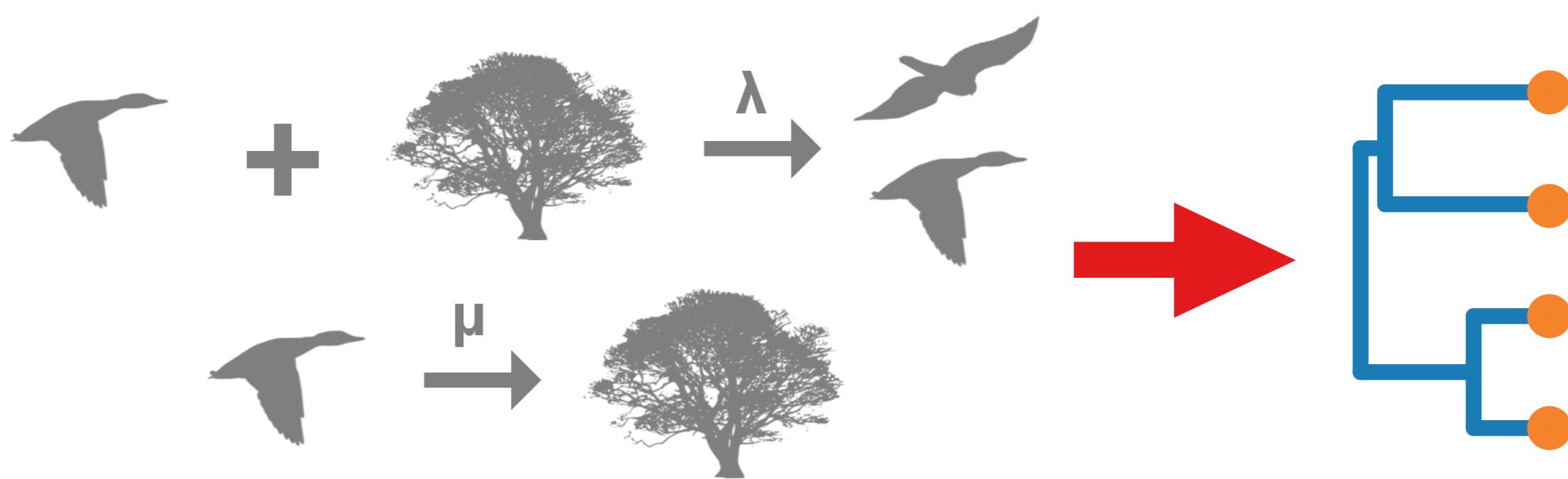
- Description of the population dynamics
 - How does the population grow over time?
 - Infectious diseases - transmission/recovery
 - Species - speciation/extinction
 - Migration events
 - Different classes of individuals
 - Host/vector dynamics
- ...etc**





Demographic model

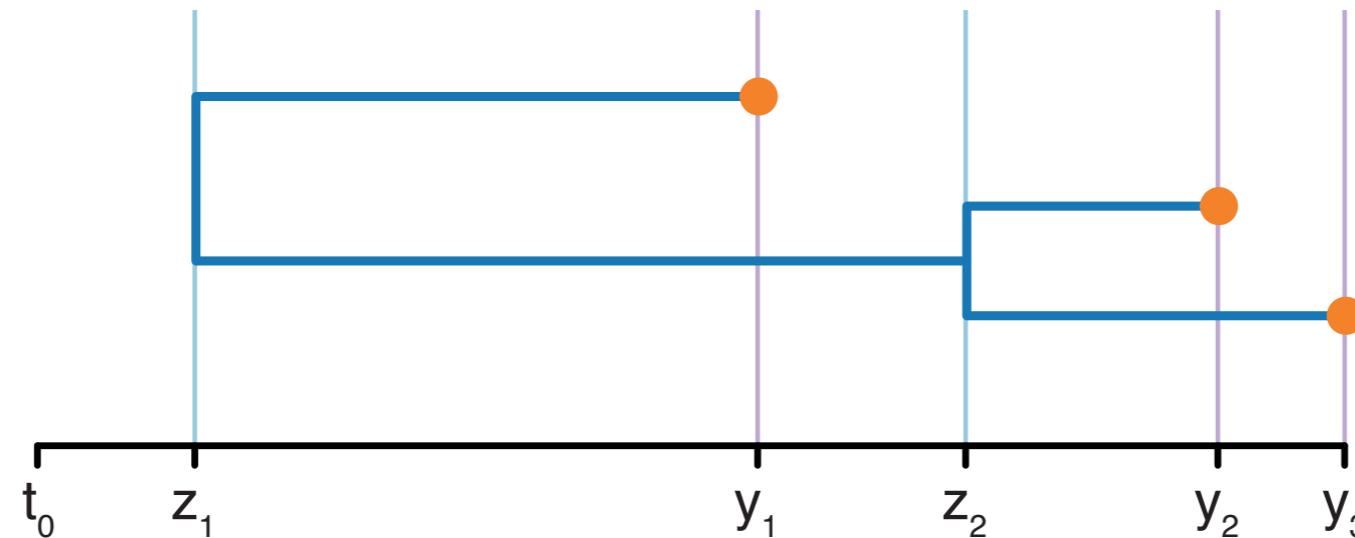
- Description of the population dynamics
 - How does the population grow over time?
 - Infectious diseases - transmission/recovery
 - Species - speciation/extinction
 - Migration events
 - Different classes of individuals
 - Host/vector dynamics
- ...etc**





Types of demographic models

Birth-death →

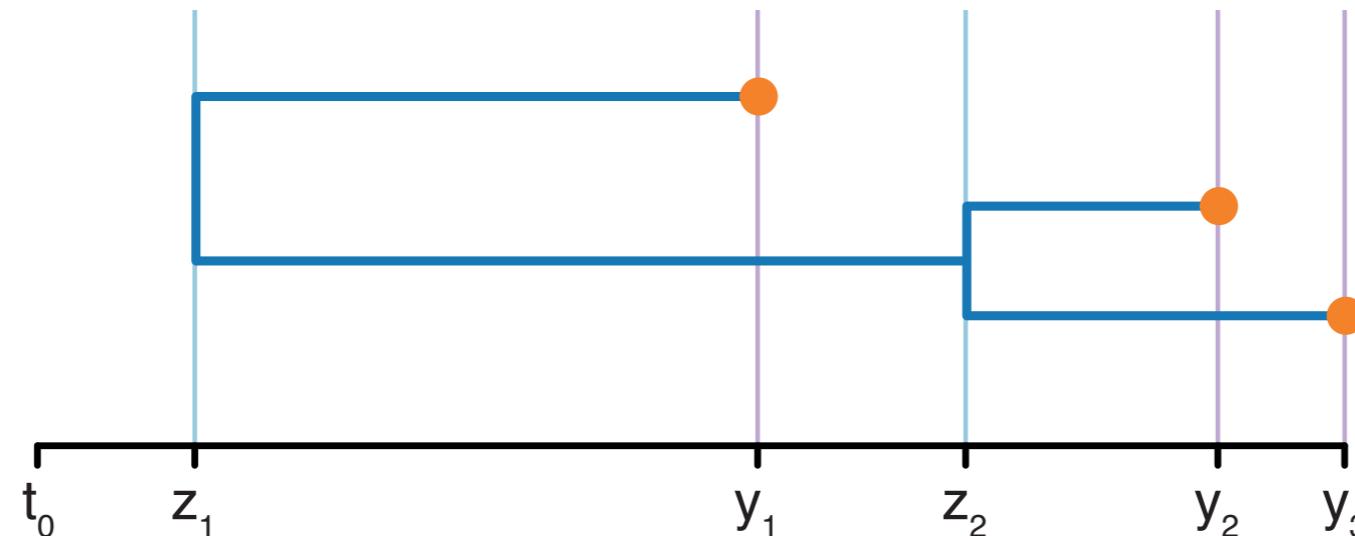


← **Coalescent**



Types of demographic models

Birth-death →



← Coalescent

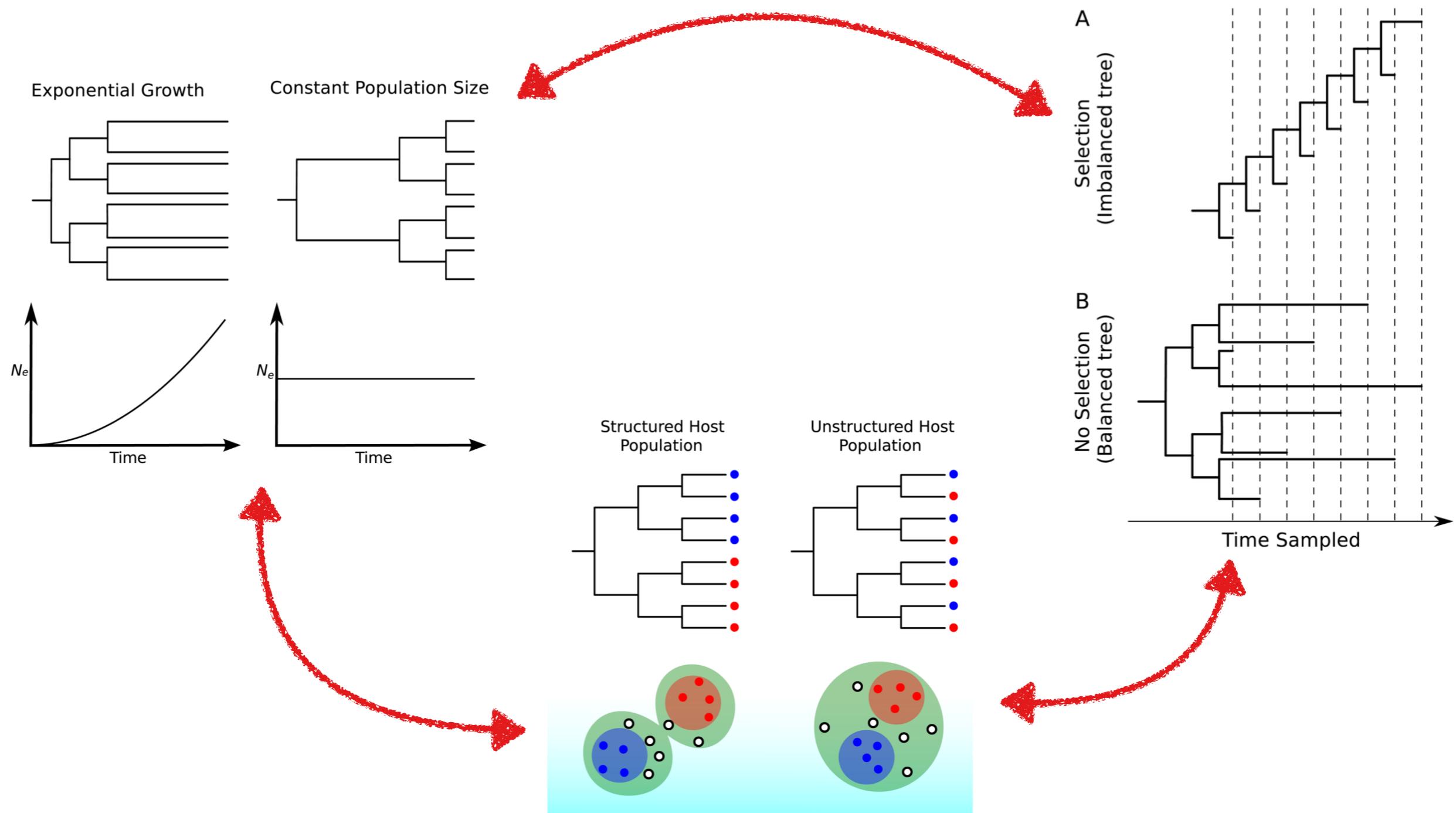
Birth-death models

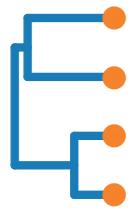
- Speciation/infection rate
- Extinction/recovery rate
- Explicitly models sampling

Coalescent

- Effective population size
- Conditions on sampling a small random sample

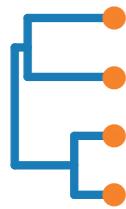
Different population dynamics generate different trees





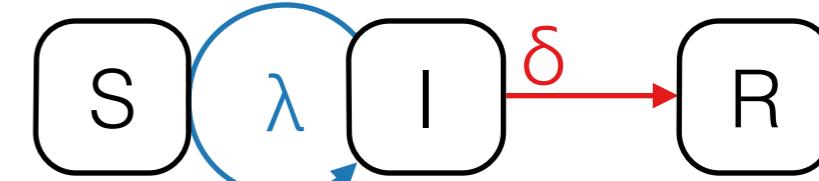
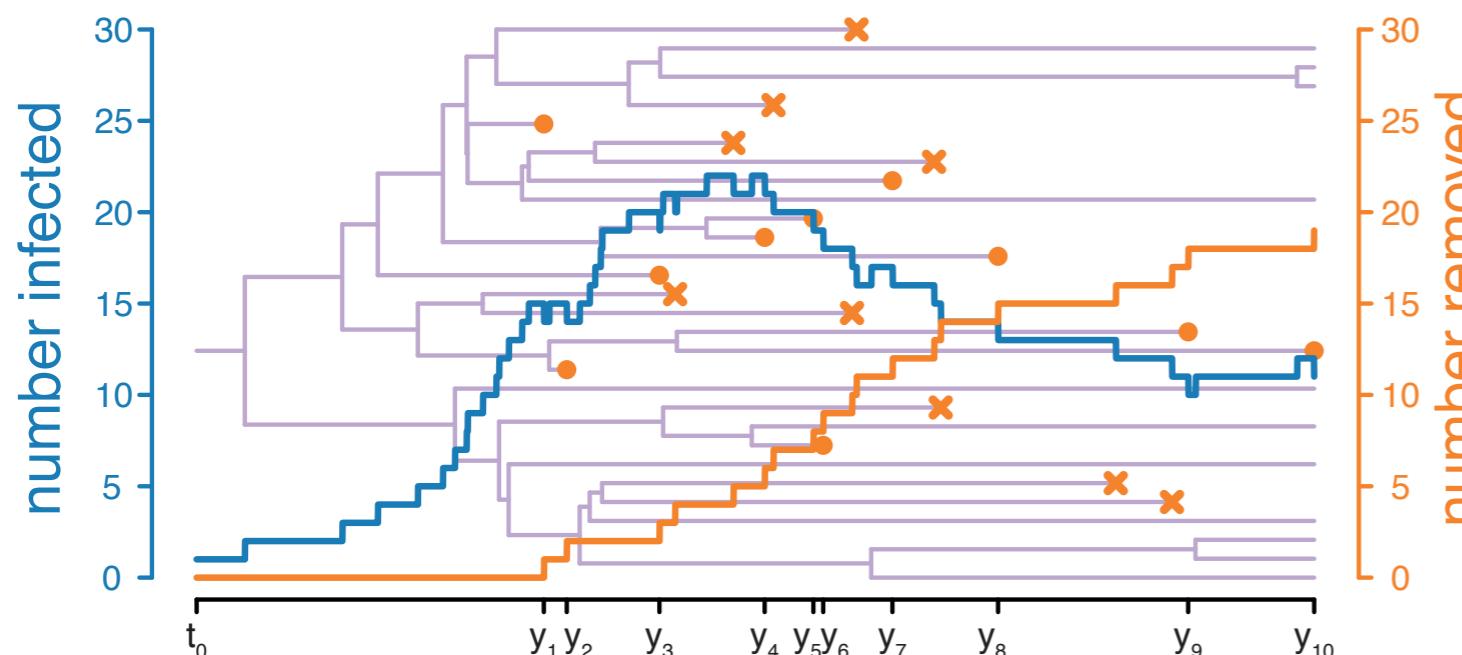
The genealogy

- What are the ancestral relationships between the sequences in our dataset?
- Only the relationships between the **sampled** sequences!

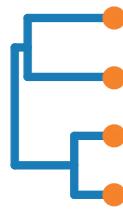


The genealogy

- What are the ancestral relationships between the sequences in our dataset?
- Only the relationships between the **sampled** sequences!

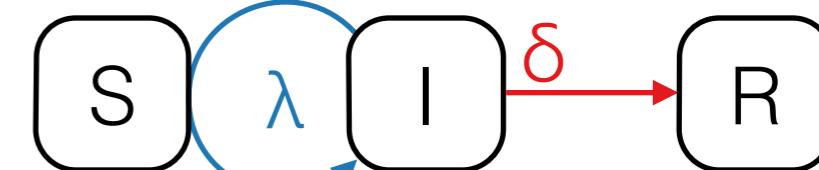
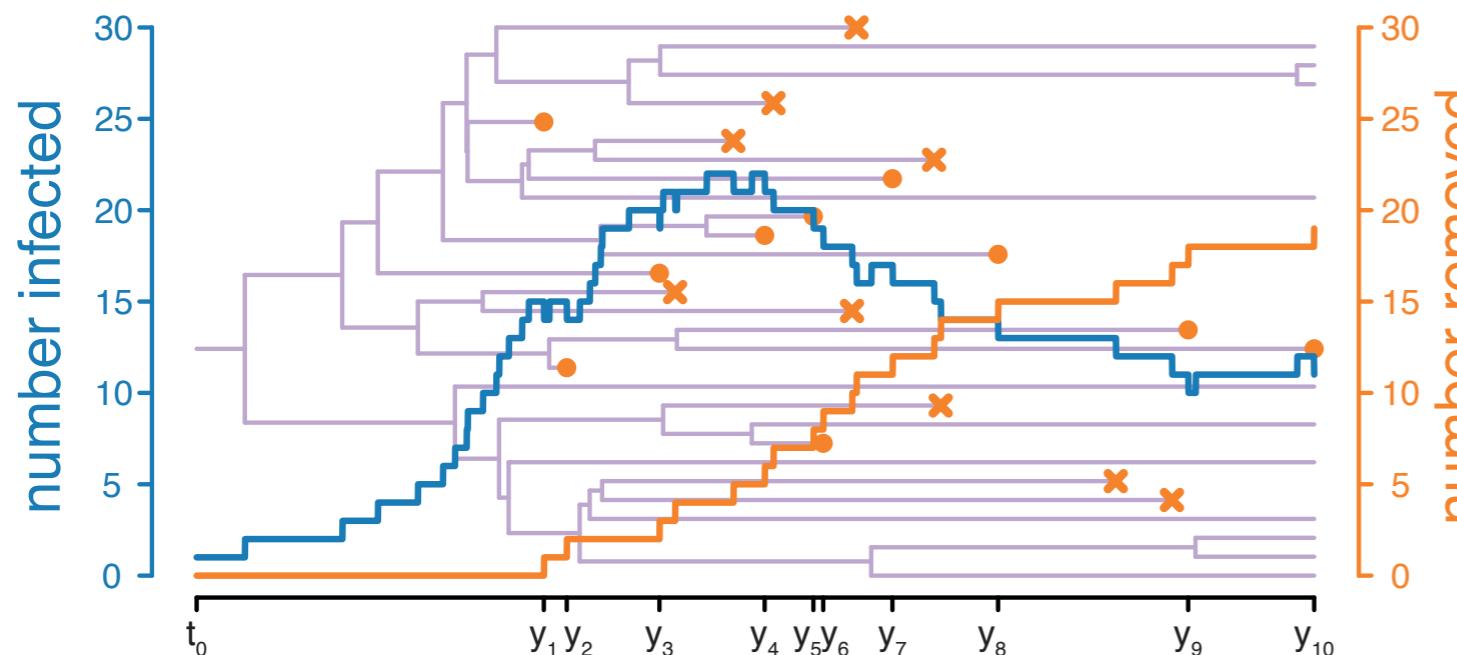


**Full tree
generated
by an SIR model**

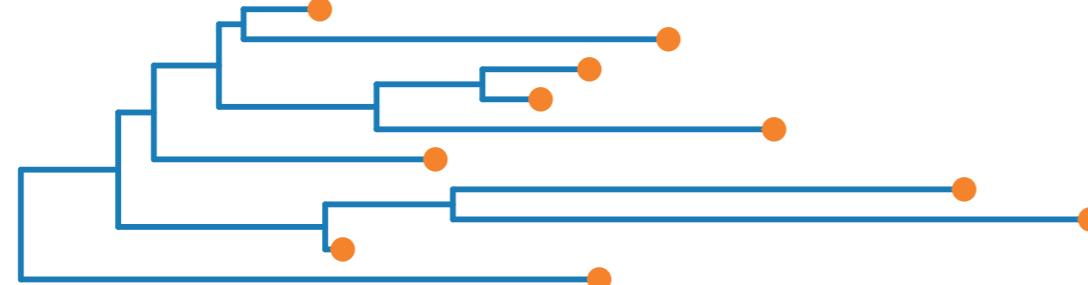


The genealogy

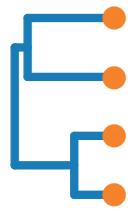
- What are the ancestral relationships between the sequences in our dataset?
- Only the relationships between the **sampled** sequences!



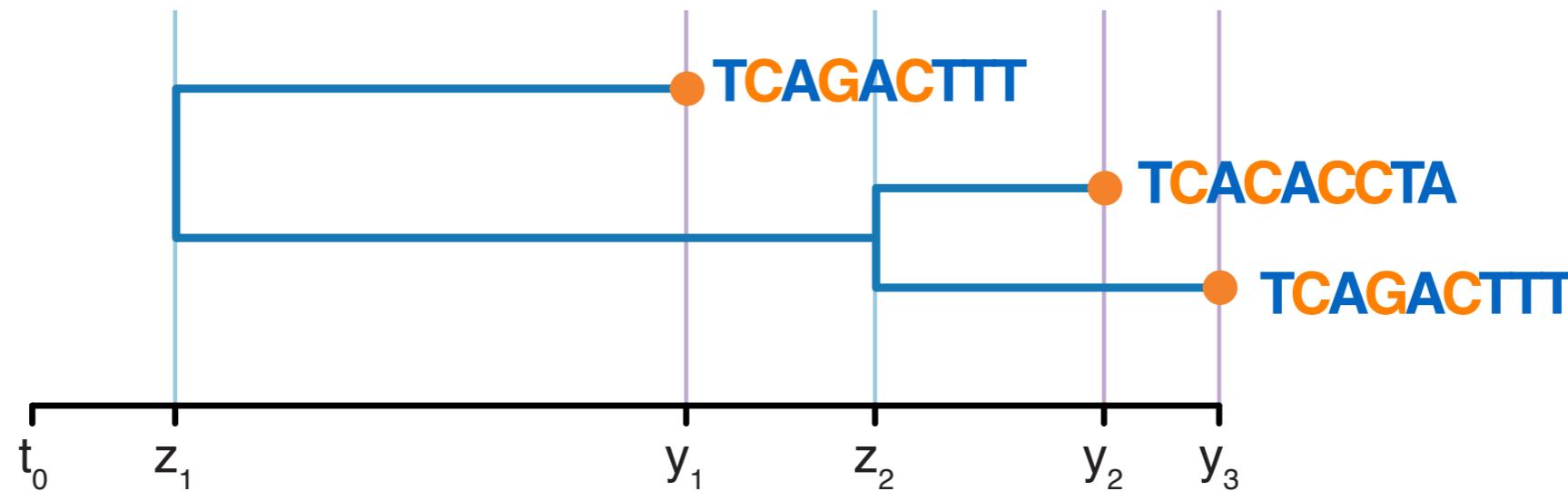
**Full tree
generated
by an SIR model**

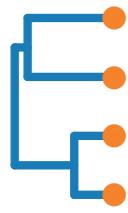


Sampled tree

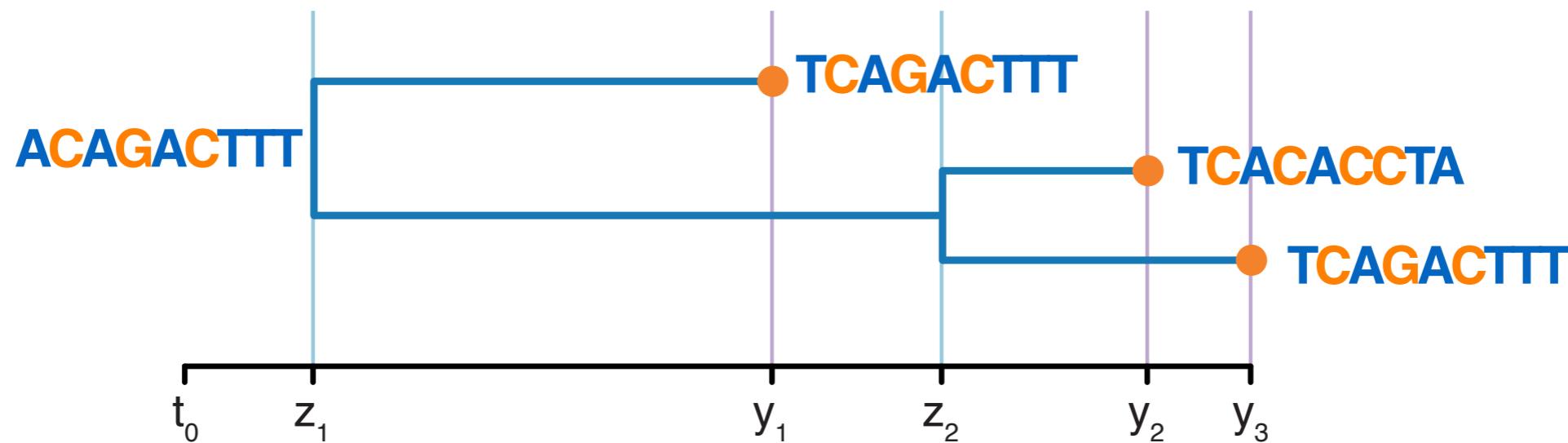


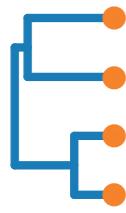
The genealogy



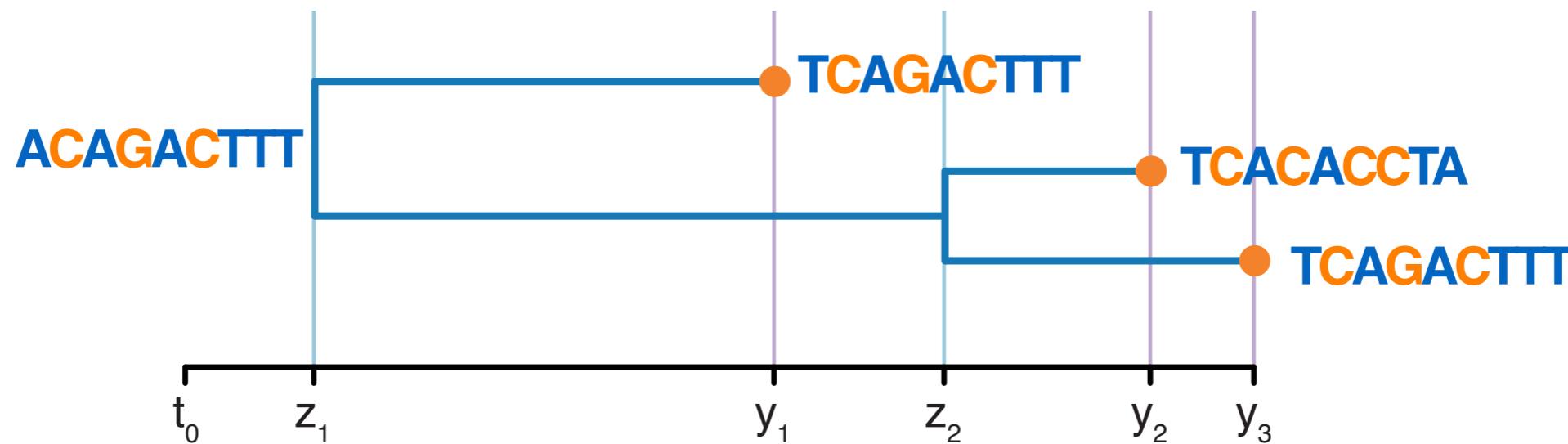


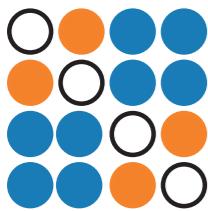
The genealogy





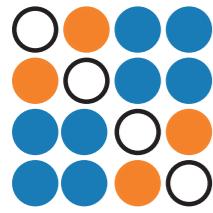
The genealogy



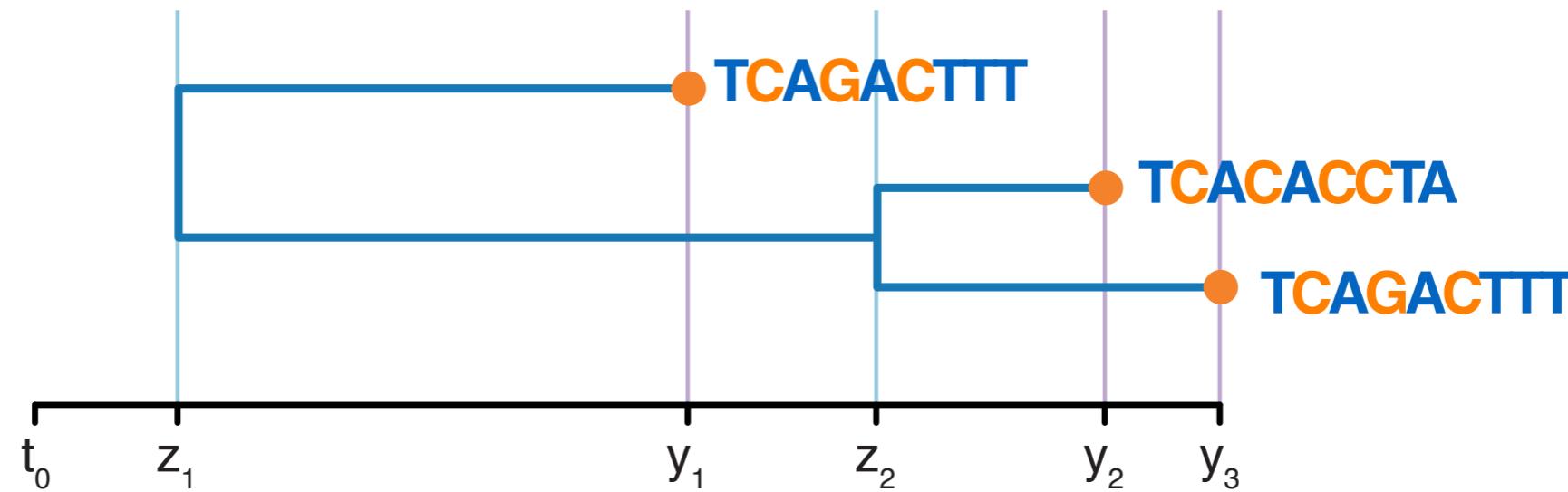


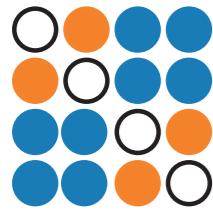
Site model

- How does one sequence change into another?
- Basic substitution model
 - Relative rate of substitution from one nucleotide to another
 - Or amino acids
 - Or codons
- Site variation
- Invariant sites
- Separate site model for each data partition
 - Multiple genes/loci
 - Model 1st, 2nd, 3rd codon positions differently

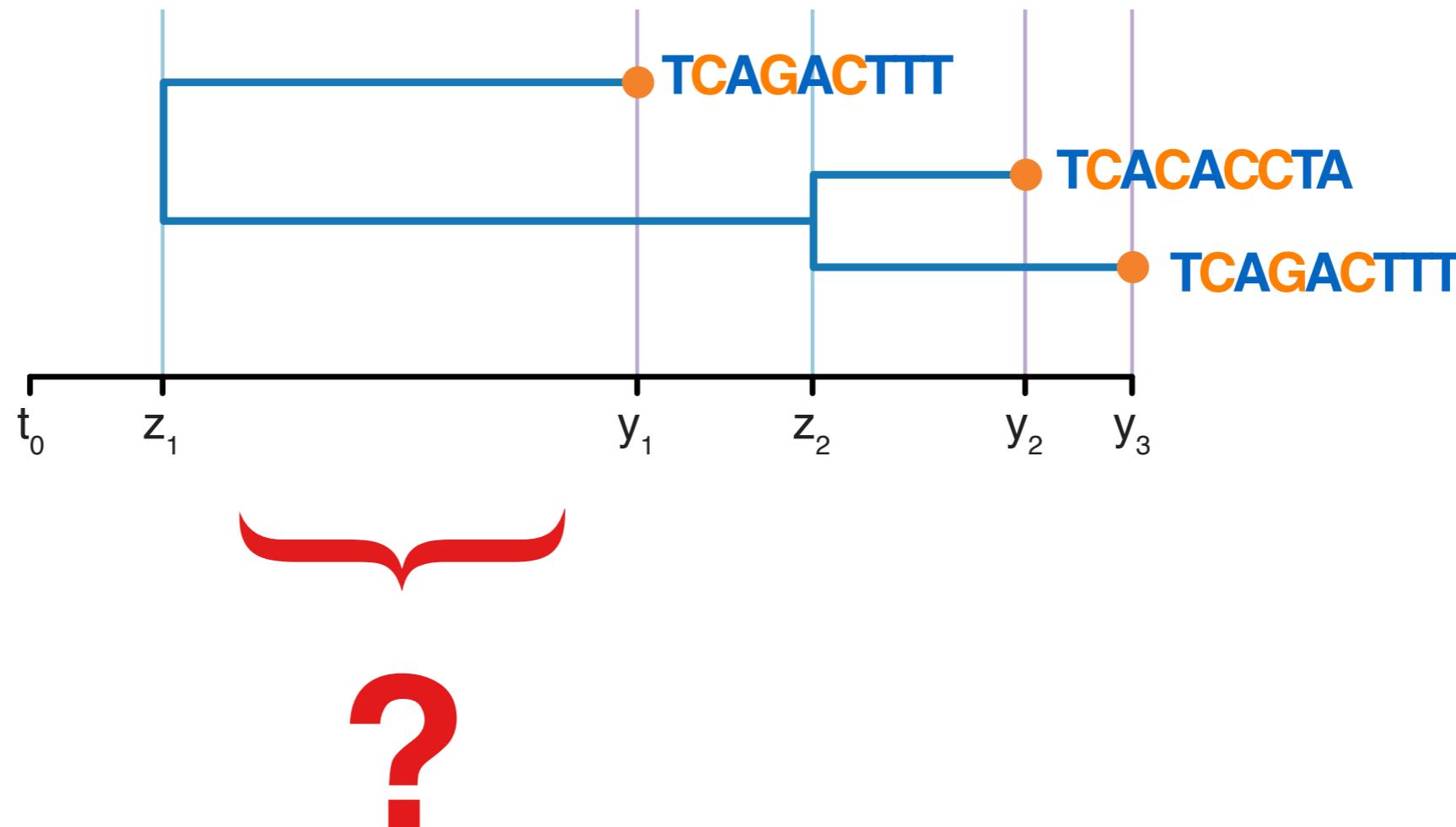


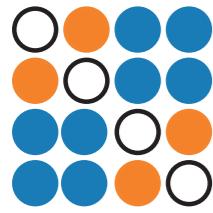
Site model



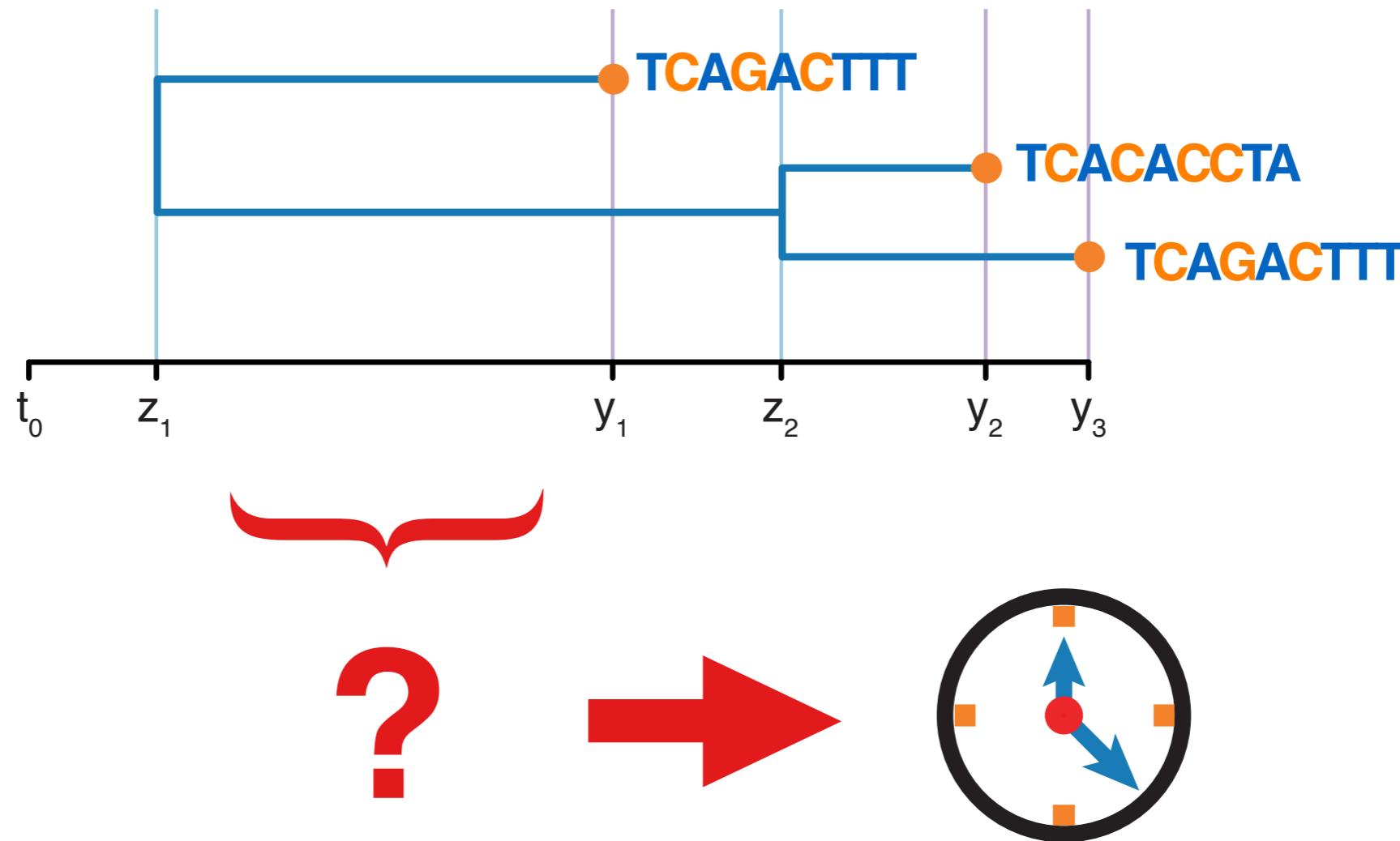


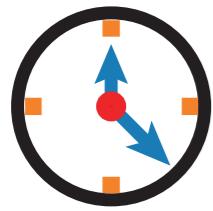
Site model





Site model





Molecular clock

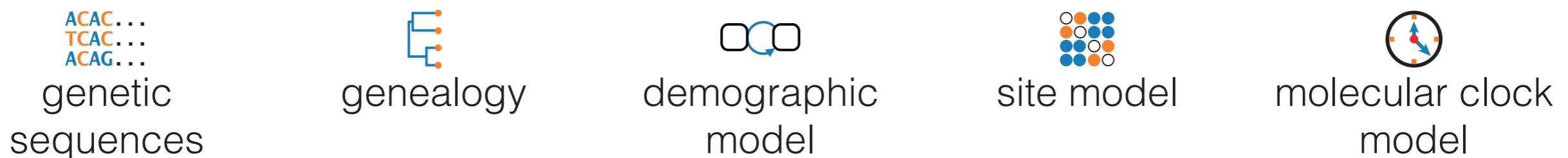
- How long does it take for substitutions to appear?
- Scales branch lengths to calendar time
- Different branches may have different clock rates

Putting it all together



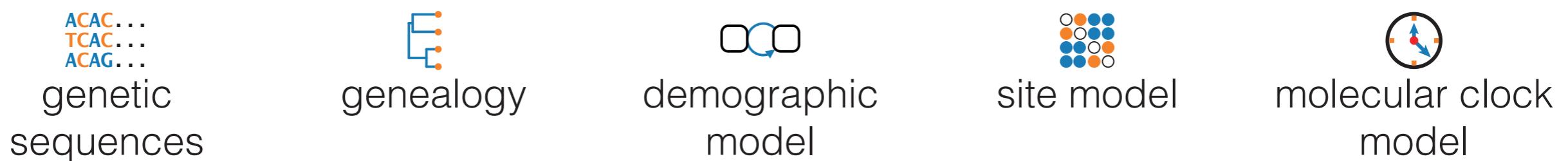
$$P(\text{model} \mid \text{data}) = \frac{P(\text{data} \mid \text{model})P(\text{model})}{P(\text{data})}$$

Putting it all together



$$P(E \circlearrowleft \text{grid} \text{ circle} | \text{ACAC...}) = \frac{P(\text{ACAC...} | E \circlearrowleft \text{grid} \text{ circle}) P(E \circlearrowleft \text{grid} \text{ circle})}{P(\text{ACAC...})}$$

Putting it all together

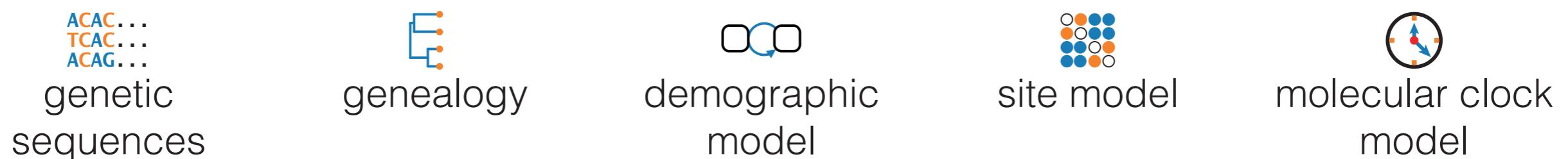


$$P(E \circ \circ \circ \bullet \bullet \bullet | ACAC \dots) = \frac{P(ACAC \dots | E \circ \circ \circ \bullet \bullet \bullet) P(E \circ \circ \circ \bullet \bullet \bullet)}{P(ACAC \dots)}$$

Assume independence

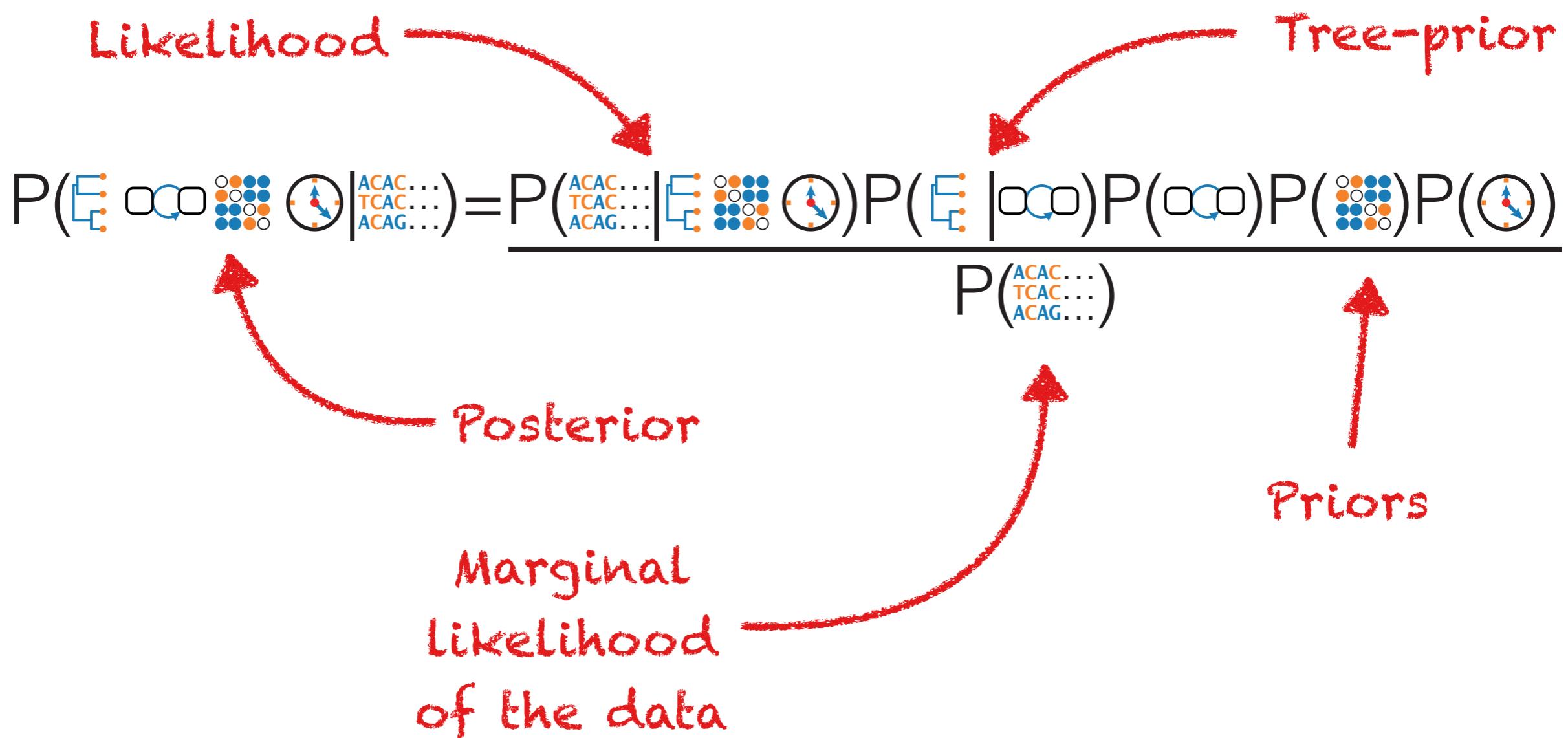
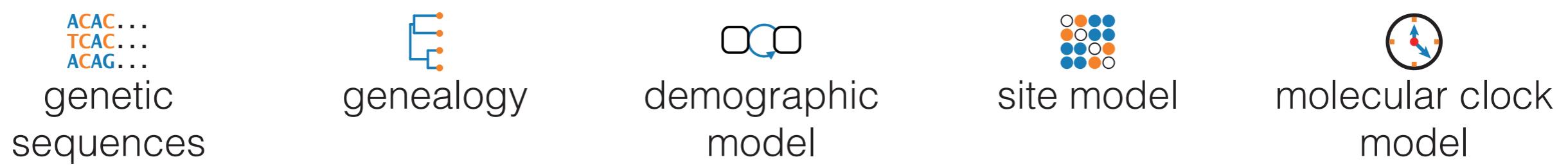
$$P(E \circ \circ \circ \bullet \bullet \bullet | \circ \circ \circ) = P(E | \circ \circ \circ) P(\circ \circ \circ) P(\bullet \bullet \bullet) P(\circ \circ \circ)$$

Posterior distribution in BEAST2



$$P(\text{genealogy} \cap \text{demographic model} \cap \text{site model} \cap \text{molecular clock model} | \text{genetic sequences}) = \frac{P(\text{genetic sequences} | \text{genealogy}, \text{demographic model}, \text{site model}, \text{molecular clock model}) P(\text{genealogy}) P(\text{demographic model}) P(\text{site model}) P(\text{molecular clock model})}{P(\text{genetic sequences})}$$

Posterior distribution in BEAST2



Posterior

ACAC...
TCAC...
ACAG...
genetic sequences

$P(E|O)$

Likelihood



MCMC inference

- We want to calculate the posterior
- But we cannot easily calculate the marginal likelihood
→ use MCMC!



molecular clock model

e-prior

$P(O|P)$

ors

Monte Carlo algorithm

- Randomized algorithm
- Deterministic runtime
(it **will** finish)
- Output may **not** be correct
(with some small probability)

MCMC inference

- MCMC performs a random walk on the posterior, preferentially sampling high-density areas
- Only need to compare which posterior density is higher
- So we only need the ratio of posteriors and the marginal likelihoods cancel out

MCMC inference

- MCMC performs a random walk on the posterior, preferentially sampling high-density areas
- Only need to compare which posterior density is higher
- So we only need the ratio of posteriors and the marginal likelihoods cancel out

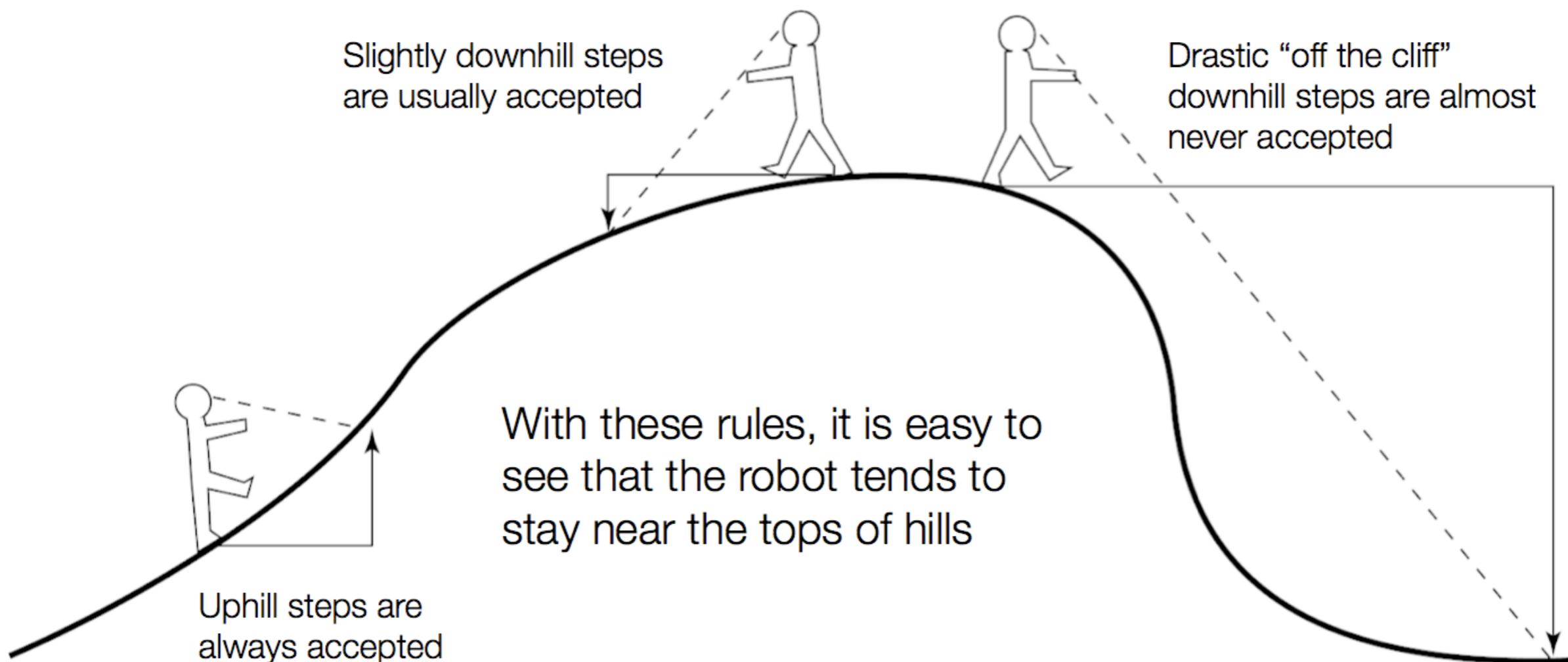
$$\frac{P(\text{model}_1 \mid \text{data})}{P(\text{model}_2 \mid \text{data})} = \frac{\frac{P(\text{data} \mid \text{model}_1)P(\text{model}_1)}{P(\text{data})}}{\frac{P(\text{data} \mid \text{model}_2)P(\text{model}_2)}{P(\text{data})}}$$

MCMC inference

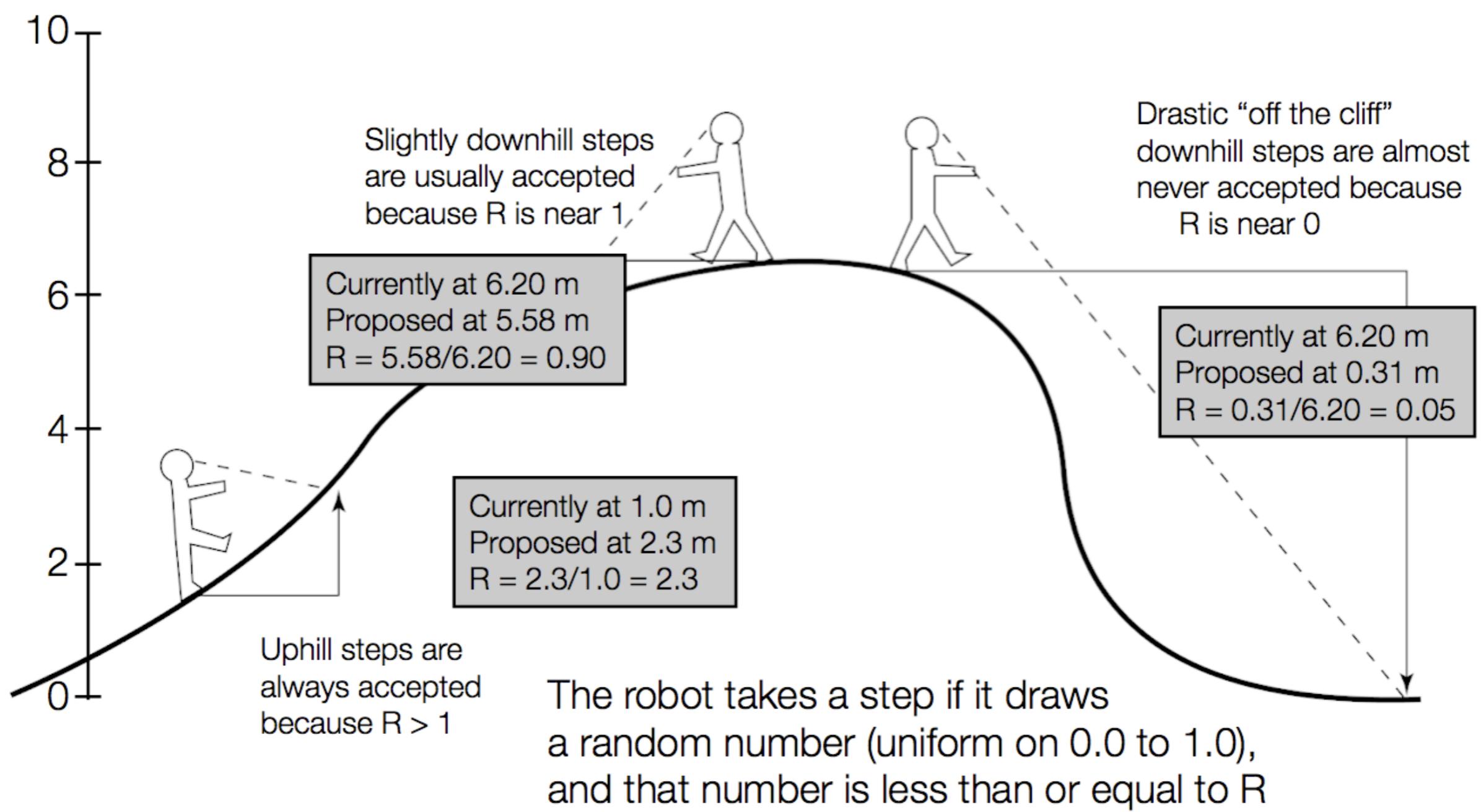
- MCMC performs a random walk on the posterior, preferentially sampling high-density areas
- Only need to compare which posterior density is higher
- So we only need the ratio of posteriors and the marginal likelihoods cancel out

$$\frac{P(\text{model}_1 \mid \text{data})}{P(\text{model}_2 \mid \text{data})} = \frac{\frac{P(\text{data} \mid \text{model}_1)P(\text{model}_1)}{P(\cancel{\text{data}})}}{\frac{P(\text{data} \mid \text{model}_2)P(\text{model}_2)}{P(\cancel{\text{data}})}}$$

MCMC robot (courtesy of Paul Lewis)

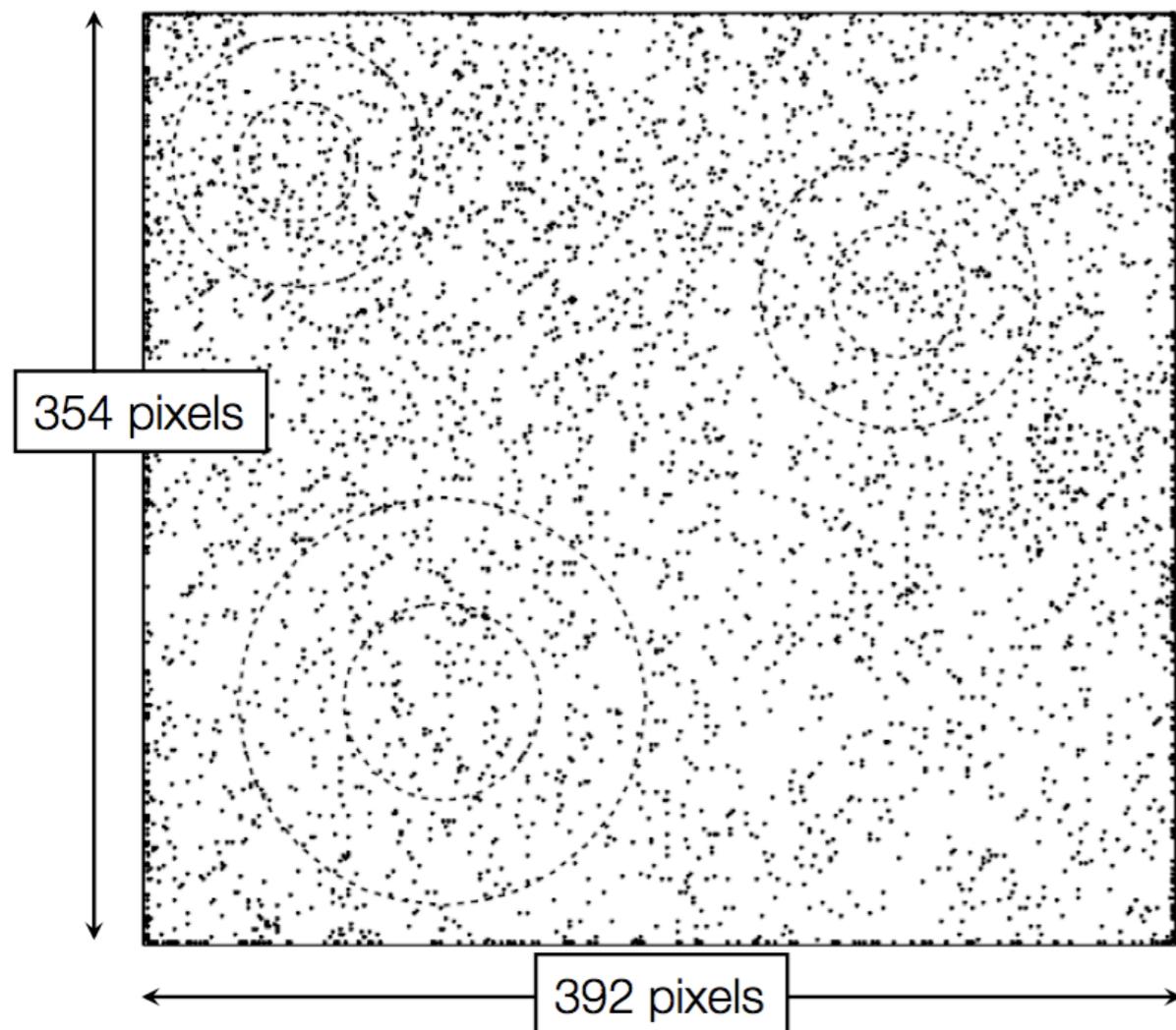


MCMC robot (courtesy of Paul Lewis)



(R is the ratio between the posterior densities)

Pure random walk (courtesy of Paul Lewis)



5000 steps by the robot

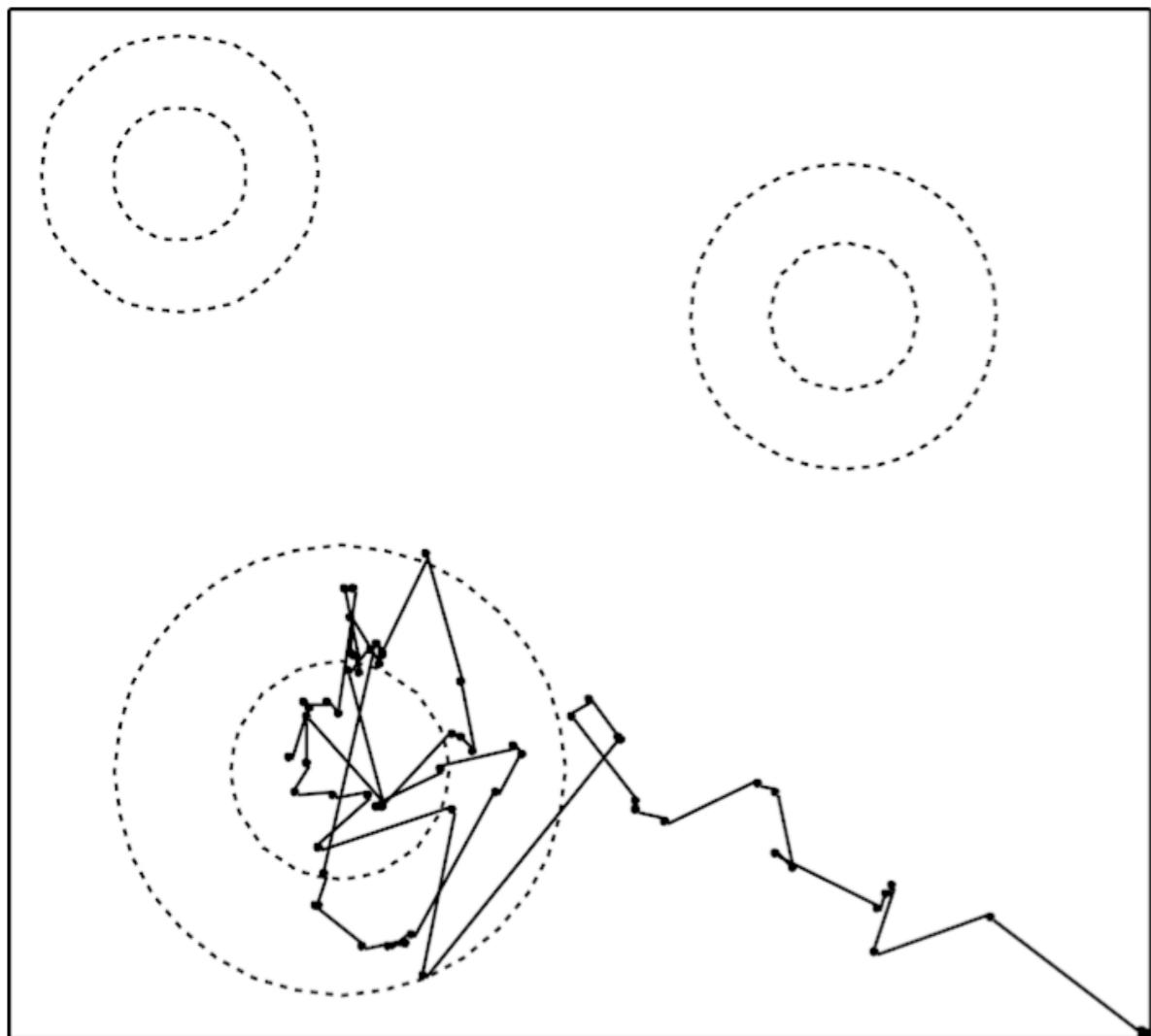
Random walk

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

Target distribution

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

Burn in (courtesy of Paul Lewis)

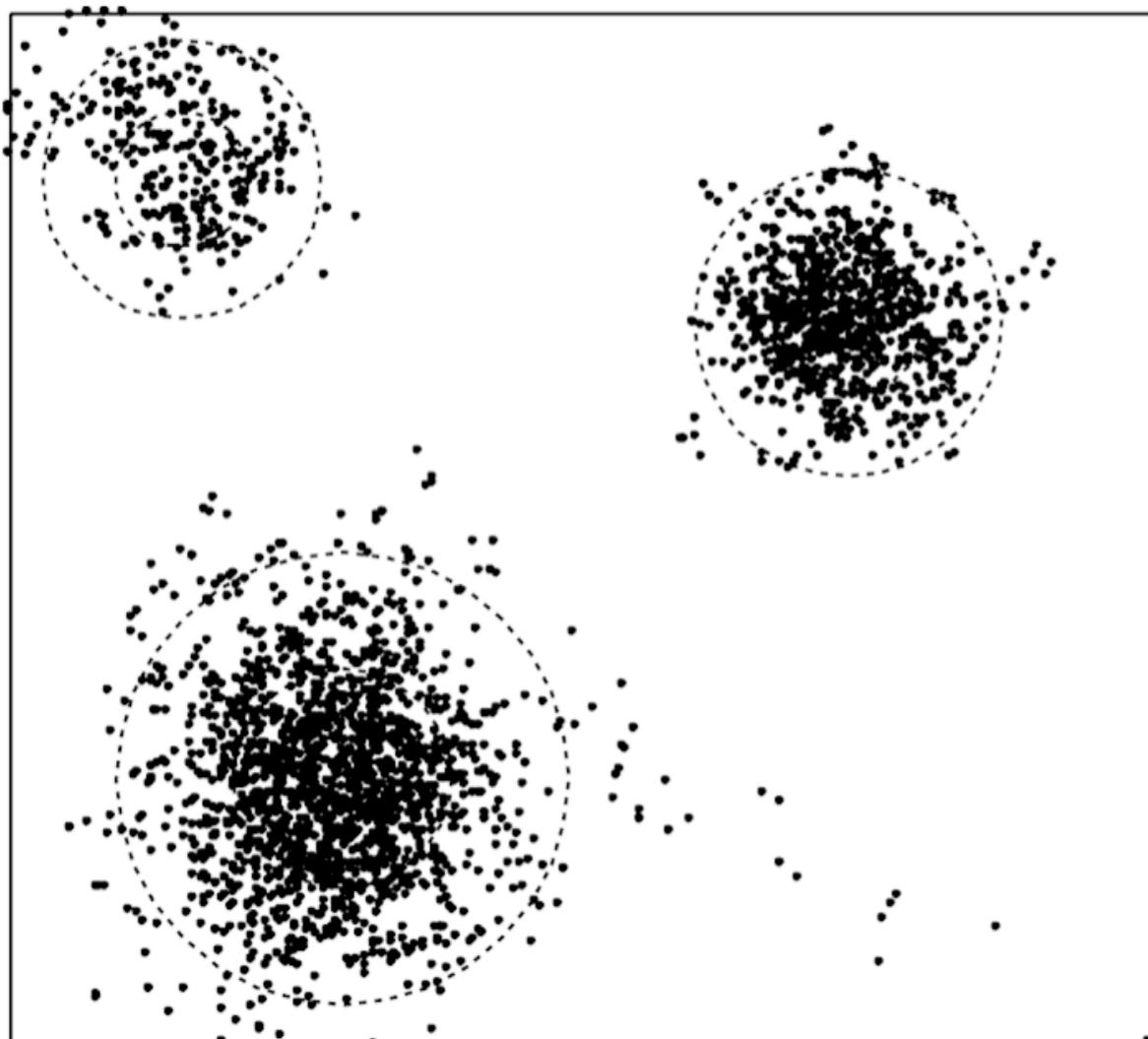


- Using MCMC rules
- Quickly finds one of the 3 hills
- First few steps are not representative of the distribution

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

Summary: MCMC inference

Target distribution

- This is the posterior in BEAST2: $P(\text{EvoTree} \mid \text{Data})$

Proposal distribution

- Used to decide where to step to next
- The choice only affects the efficiency of the algorithm
- Metropolis algorithm: symmetric proposals
- Metropolis-Hastings: asymmetric proposals

Summary: MCMC inference

Target distribution

- This is the posterior in BEAST2: $P(E \mid \text{data})$

Proposal distribution

- Used to decide where to step to next
 - The choice only affects the efficiency of the algorithm
 - Metropolis algorithm: symmetric proposals
 - Metropolis-Hastings: asymmetric proposals

In BEAST and BEAST2 operators are used to propose the next step

Marginal distributions

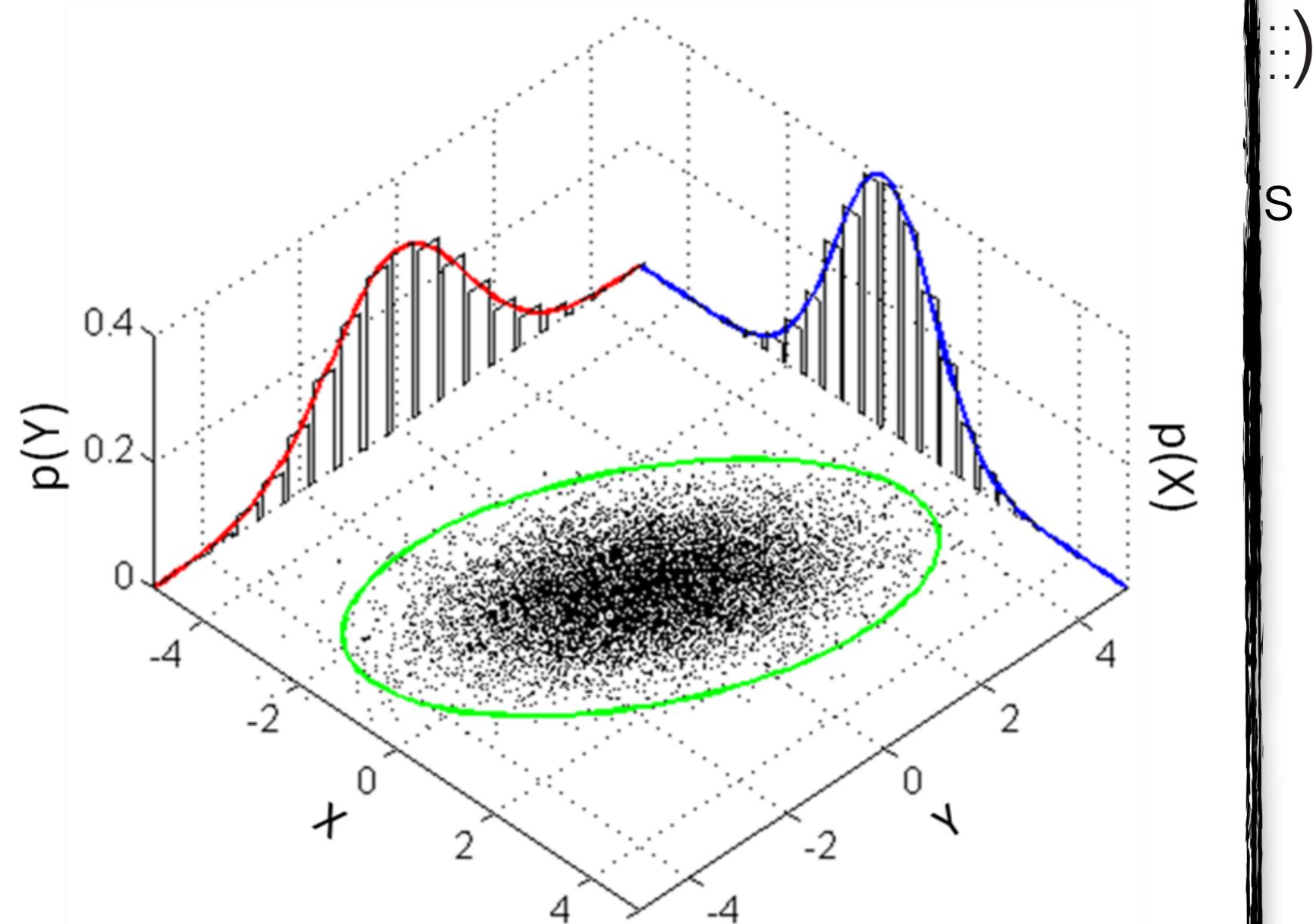
- We only have the joint posterior: $P(\text{E} \text{ } \text{ } \text{ } \text{ } \text{ } | \text{ACAC}, \text{TCAC}, \text{ACAG}, \dots)$
- But we want distributions for each of the parameters we are interested in → marginalize

$$P(\phi) = \int_{\theta} P(\phi|\theta)P(\theta)d\theta$$

Margin

In practice

- We often
- But we
- we are



One tree to rule them all?

- As with other parameters the data may support multiple trees equally well
- Thus we should look at all trees in order to take phylogenetic uncertainty into account
- If we are only interested in the epidemiological parameters we can integrate out the tree (but this is very difficult in an ML framework)
- Naturally integrate out the tree in an MCMC framework by inferring the distribution of trees

$$P(\text{ACAC...} | \text{TCAC...} \text{ ACAG...}) = \frac{P(\text{ACAC...} | \text{TCAC...} \text{ ACAG...}) P(\text{genealogy}) P(\text{demographic model}) P(\text{site model}) P(\text{molecular clock model})}{P(\text{ACAC...} | \text{TCAC...} \text{ ACAG...})}$$

Legend:

- genetic sequences:
- genealogy:
- demographic model:
- site model:
- molecular clock model:

One tree to rule them all?

- As with other parameters the data may support multiple trees equally well
- Thus we should look at all trees in order to take

Can easily integrate out any nuisance parameters!

- What is a nuisance parameter depends on the question
- In phylogenetics we specifically want the tree, but the substitution model is a nuisance parameter



$$P(\overset{\text{ACAC}}{\text{TCAC}} \overset{\text{TCAC}}{\text{ACAG}} \dots)$$

genetic sequences

genealogy

demographic model

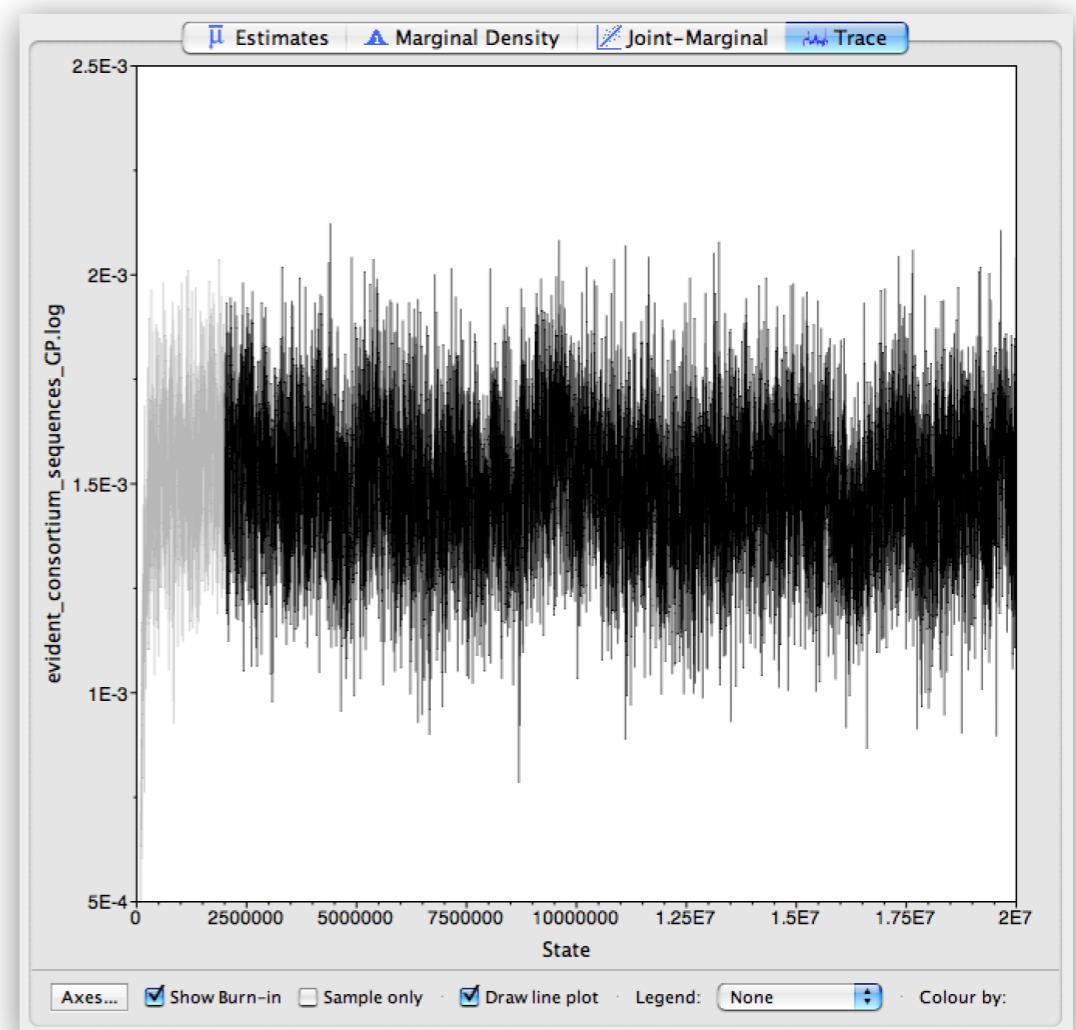
site model

molecular clock model

What we hope will happen

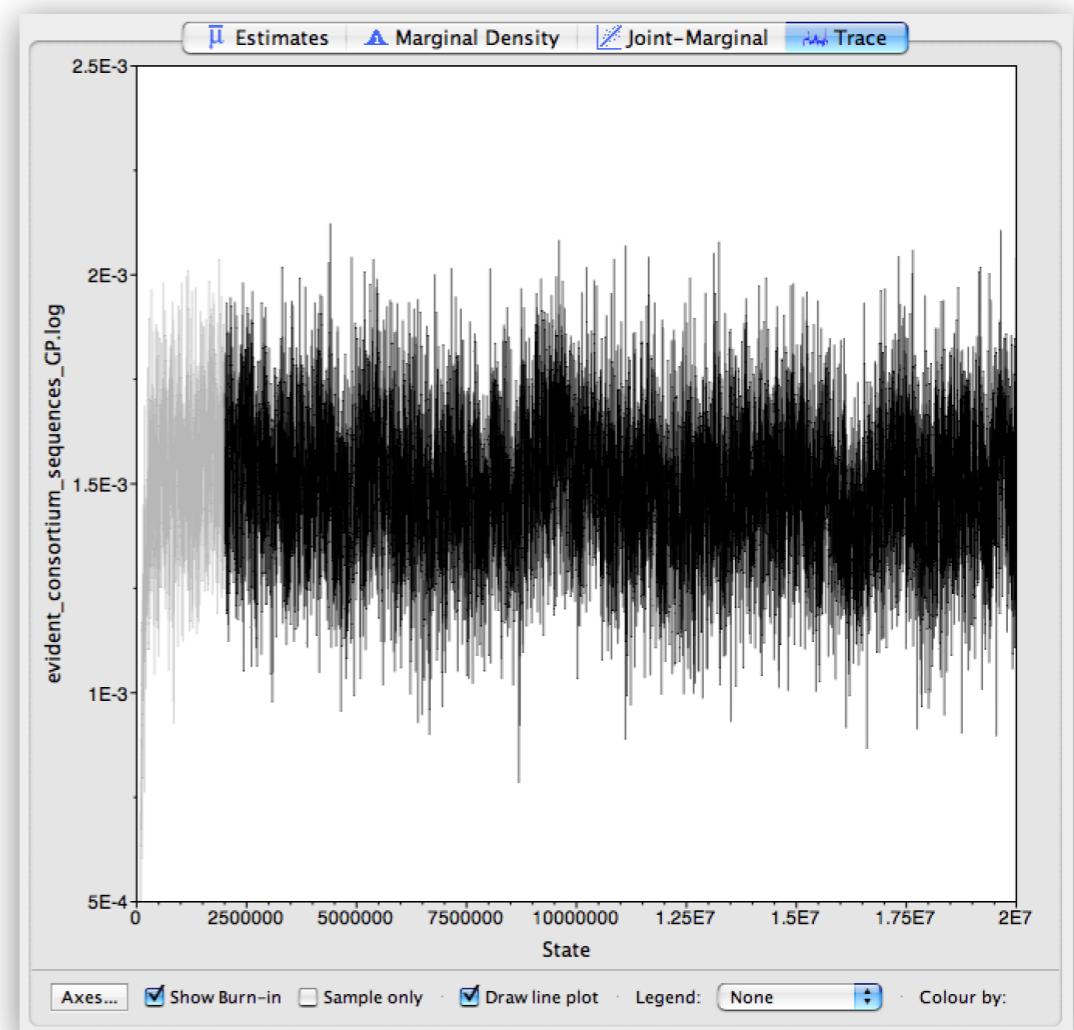
What we hope will happen

- The MCMC algorithm samples efficiently from high density areas of the posterior distribution



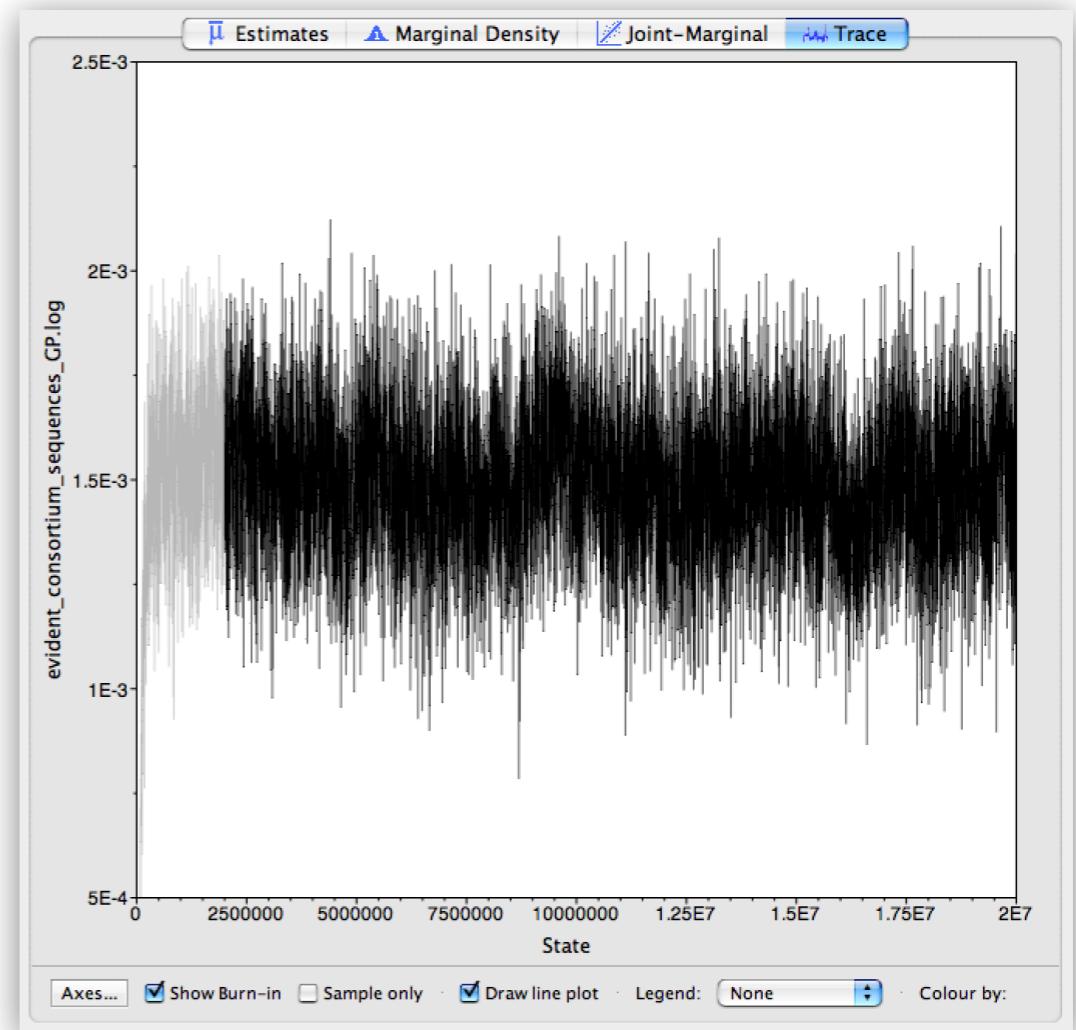
What we hope will happen

- 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



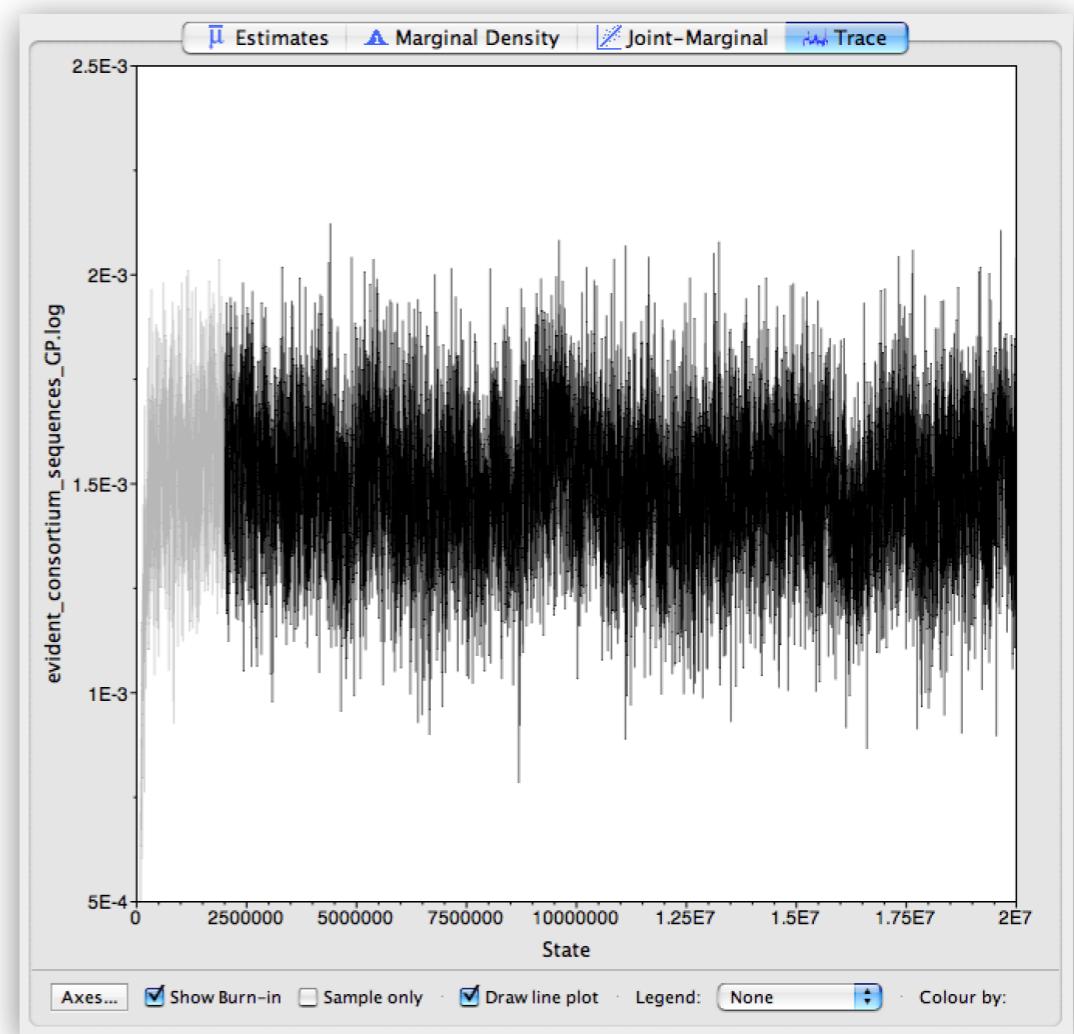
What we hope will happen

- 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
- Everything is awesome!



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
- Everything is awesome!

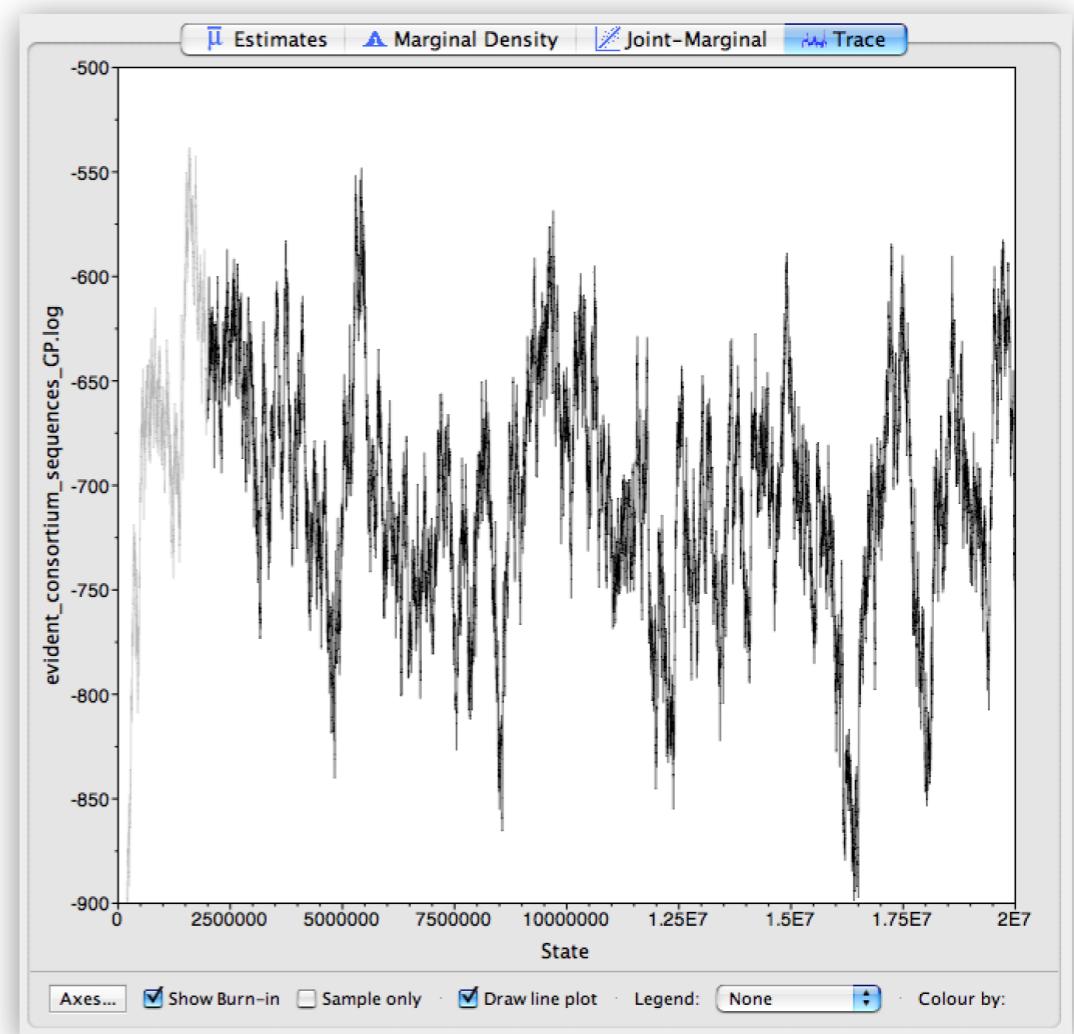


Mixing well! 😊

What often happens in practice...

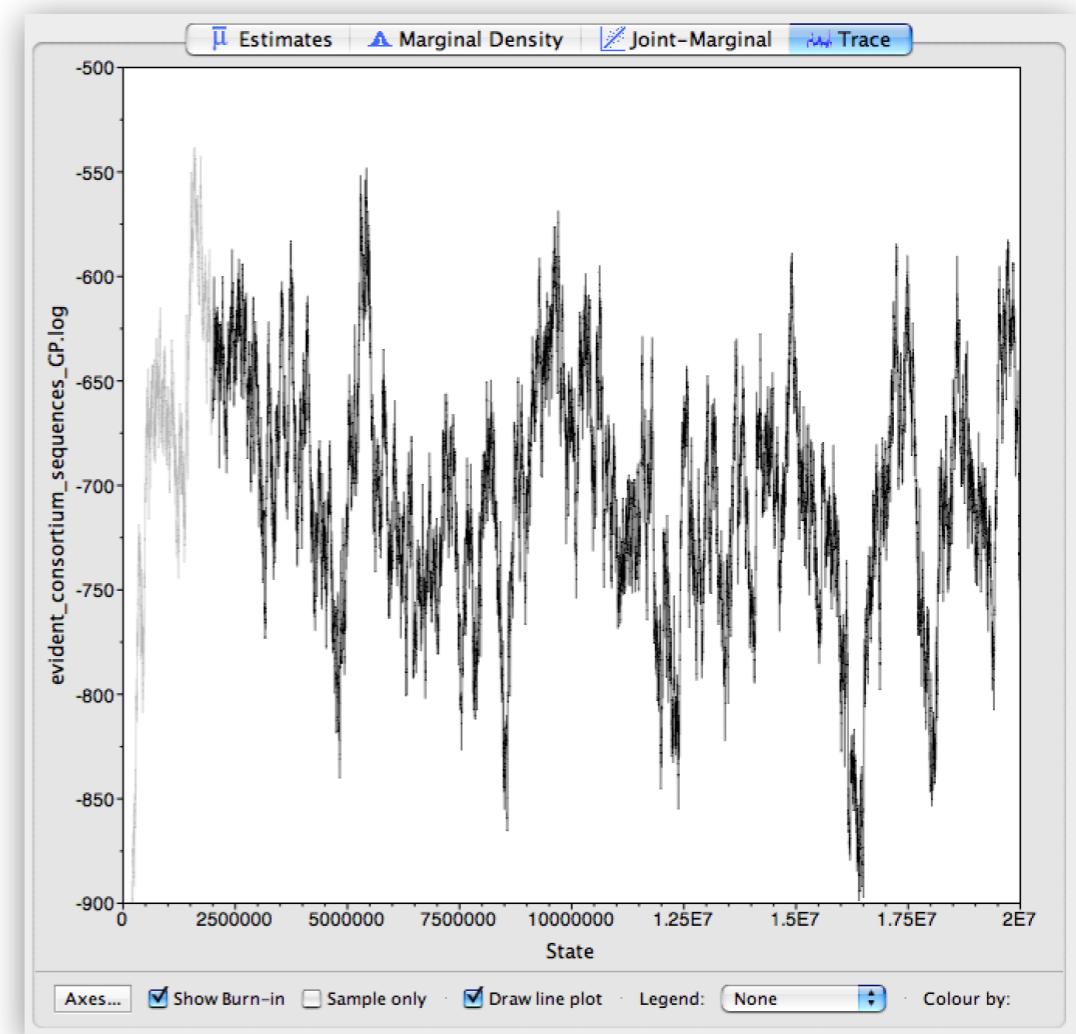
What often happens in practice...

- Not so much...



What often happens in practice...

- Not so much...



Not mixing! 😭

What often goes wrong?

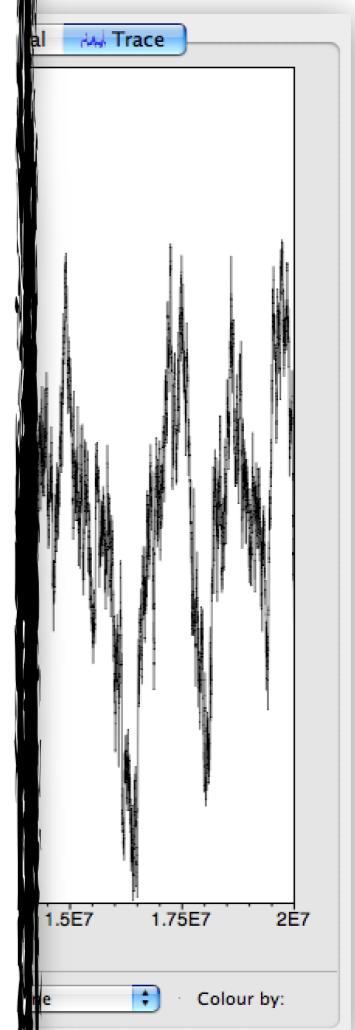
- Not so much

What went wrong?

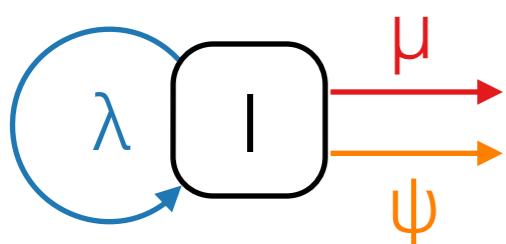
- Recall that MCMC is a Monte Carlo algorithm — it is not **guaranteed** to find the correct solution in finite time!

How can we make it work better?

- Tweak the operators to make good proposals
(increase operator efficiency)
- Fine tune specifics about the MCMC run
(sampling frequency, length etc.)
- Poor model choice?
- Poor choice of parameterization?

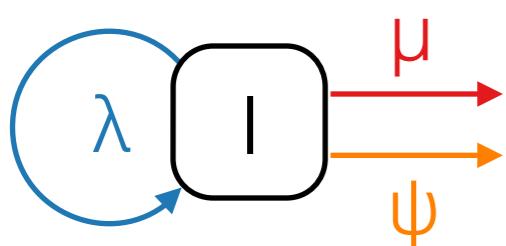


Model Parameterization: Birth-death process



- λ — Birth-rate (speciation / transmission)
- μ — Death-rate (extinction / death)
- ψ — Sampling rate

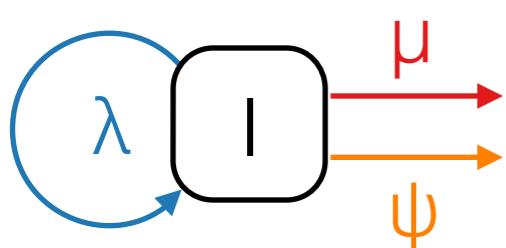
Model Parameterization: Birth-death process



- λ — Birth-rate (speciation / transmission) $\lambda \in [0, \infty)$
- μ — Death-rate (extinction / death) $\mu \in [0, \infty)$
- ψ — Sampling rate $\psi \in [0, \infty)$

Priors?

Model Parameterization: Birth-death process



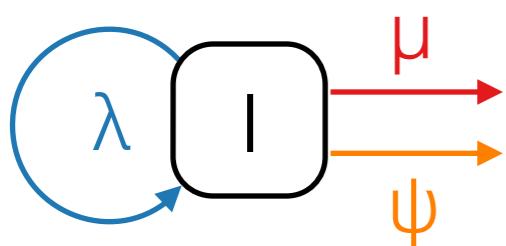
- λ — Birth-rate (speciation / transmission) $\lambda \in [0, \infty)$
- μ — Death-rate (extinction / death) $\mu \in [0, \infty)$
- ψ — Sampling rate $\psi \in [0, \infty)$

Priors?

Infectious diseases

- $R = \lambda/(\mu + \psi)$ — Reproduction number
(is it spreading or not?)
- $\delta = \mu + \psi$ — Becoming noninfectious rate
($1/\delta$ = infectious period)
- $p = \psi/(\mu + \psi)$ — Sampling proportion

Model Parameterization: Birth-death process



- λ — Birth-rate (speciation / transmission) $\lambda \in [0, \infty)$
- μ — Death-rate (extinction / death) $\mu \in [0, \infty)$
- ψ — Sampling rate $\psi \in [0, \infty)$

Priors?

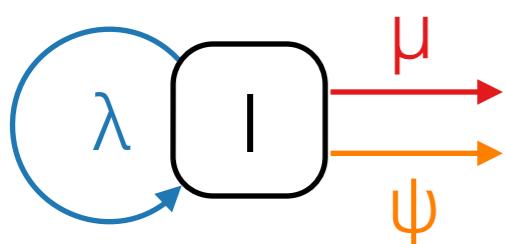
Infectious diseases

- $R = \lambda/(\mu + \psi)$ — Reproduction number
(is it spreading or not?)
- $\delta = \mu + \psi$ — Becoming noninfectious rate
($1/\delta$ = infectious period)
- $p = \psi/(\mu + \psi)$ — Sampling proportion

$$R \in [0, \infty) \quad \delta \in [0, \infty) \quad p \in [0, 1]$$

More natural priors!

Model Parameterization: Birth-death process



- λ — Birth-rate (speciation / transmission) $\lambda \in [0, \infty)$
- μ — Death-rate (extinction / death) $\mu \in [0, \infty)$
- ψ — Sampling rate $\psi \in [0, \infty)$

Priors?

Infectious diseases

- $R = \lambda/(\mu + \psi)$ — Reproduction number
(is it spreading or not?)
- $\delta = \mu + \psi$ — Becoming noninfectious rate
($1/\delta$ = infectious period)
- $p = \psi/(\mu + \psi)$ — Sampling proportion

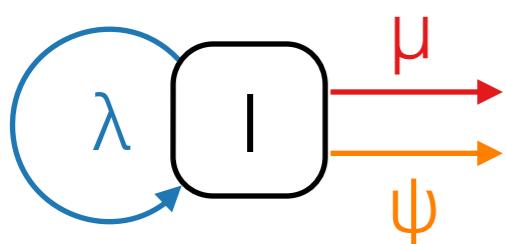
$$R \in [0, \infty) \quad \delta \in [0, \infty) \quad p \in [0, 1]$$

Speciation

- $\lambda - \mu$ — Growth rate
- μ/λ — Relative death rate
- $p = \psi/(\mu + \psi)$ — Sampling proportion

More natural priors!

Model Parameterization: Birth-death process



- λ — Birth-rate (speciation / transmission) $\lambda \in [0, \infty)$
- μ — Death-rate (extinction / death) $\mu \in [0, \infty)$
- ψ — Sampling rate $\psi \in [0, \infty)$

Priors?

Infectious diseases

- $R = \lambda/(\mu + \psi)$ — Reproduction number
(is it spreading or not?)
- $\delta = \mu + \psi$ — Becoming noninfectious rate
($1/\delta$ = infectious period)
- $p = \psi/(\mu + \psi)$ — Sampling proportion

$$R \in [0, \infty) \quad \delta \in [0, \infty) \quad p \in [0, 1]$$

More natural priors!

Speciation

- $\lambda - \mu$ — Growth rate
- μ/λ — Relative death rate
- $p = \psi/(\mu + \psi)$ — Sampling proportion

$$(\lambda - \mu) \in [0, \infty) \quad \mu/\lambda \in [0, 1] \quad p \in [0, 1]$$

Smaller parameter space!

Summary

- Model consists of a demographic model, a site model and a molecular clock model
- Combine priors for model parameters with the likelihood to get the posterior
- Using MCMC we can efficiently approximate the joint posterior of the model
- By marginalising we get posterior distributions of parameters of interest while integrating out uncertainty in nuisance parameters
- Need a good set of operators to propose new steps
- Efficient inference may need some fine-tuning!
(or it may not converge to the posterior distribution in a reasonable time)