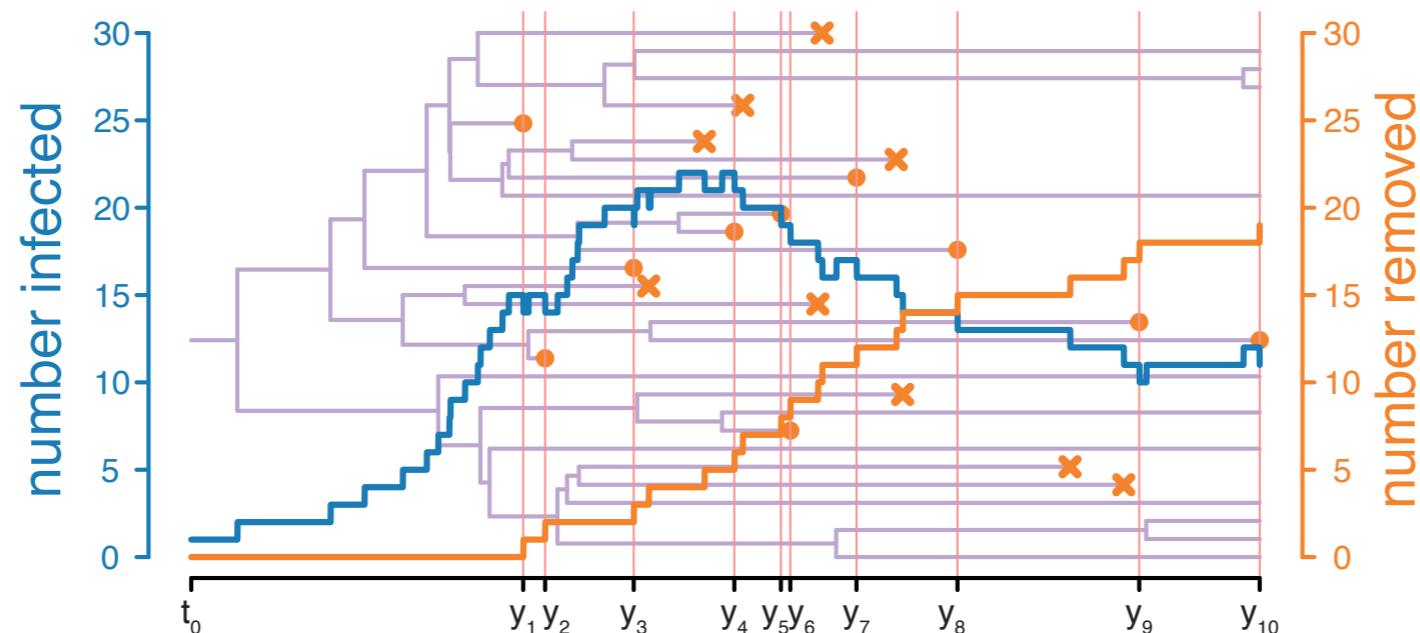


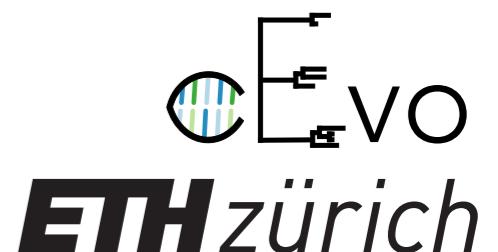


# Inferring epidemiological parameters from phylogenies with BEAST2



Swiss Institute of  
Bioinformatics

Veronika Bošková, Louis du Plessis, Joëlle Barido-Sottani  
Computational Evolution ([www.bsse.ch/cevo](http://www.bsse.ch/cevo))  
Department of Biosystems Science and Engineering





# Inferring epidemiological parameters from phylogenies with BEAST2

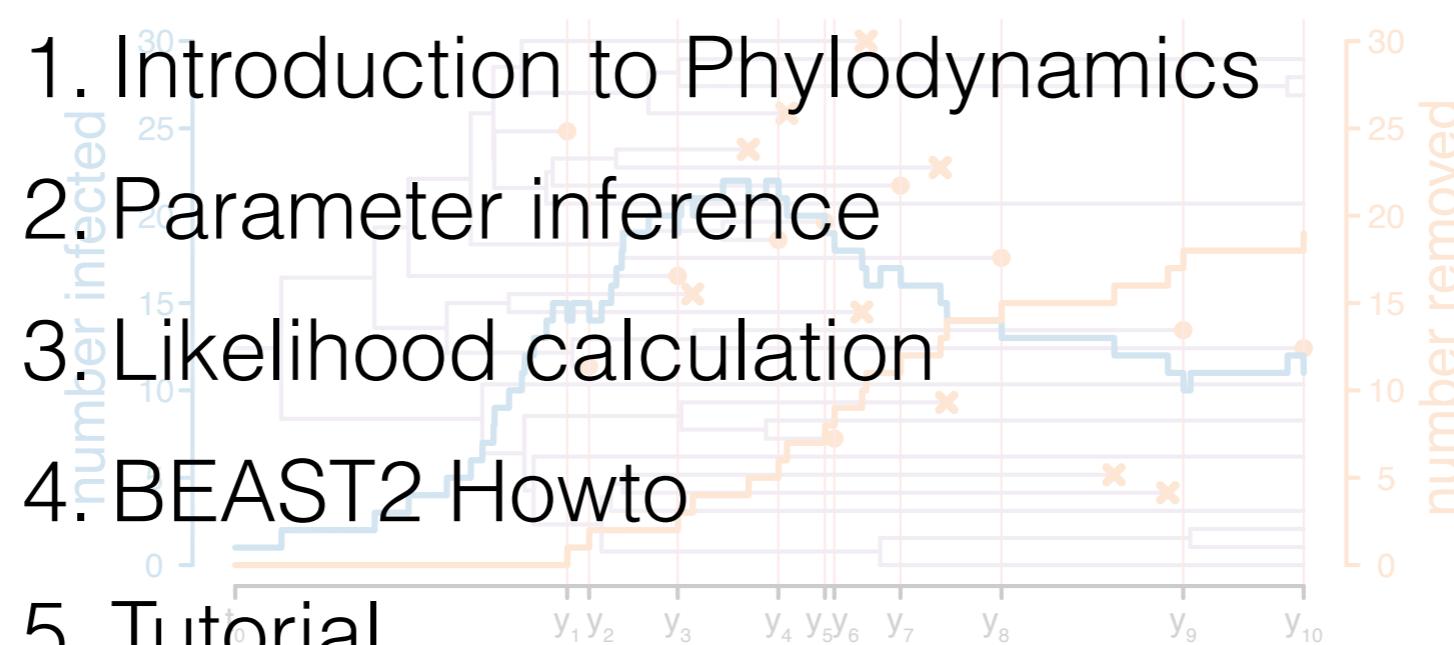
1. Introduction to Phylodynamics

2. Parameter inference

3. Likelihood calculation

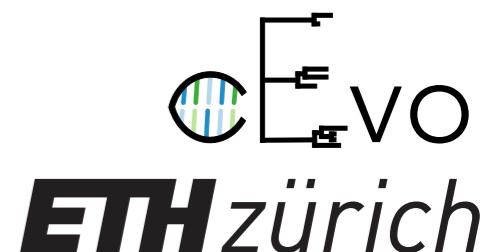
4. BEAST2 Howto

5. Tutorial



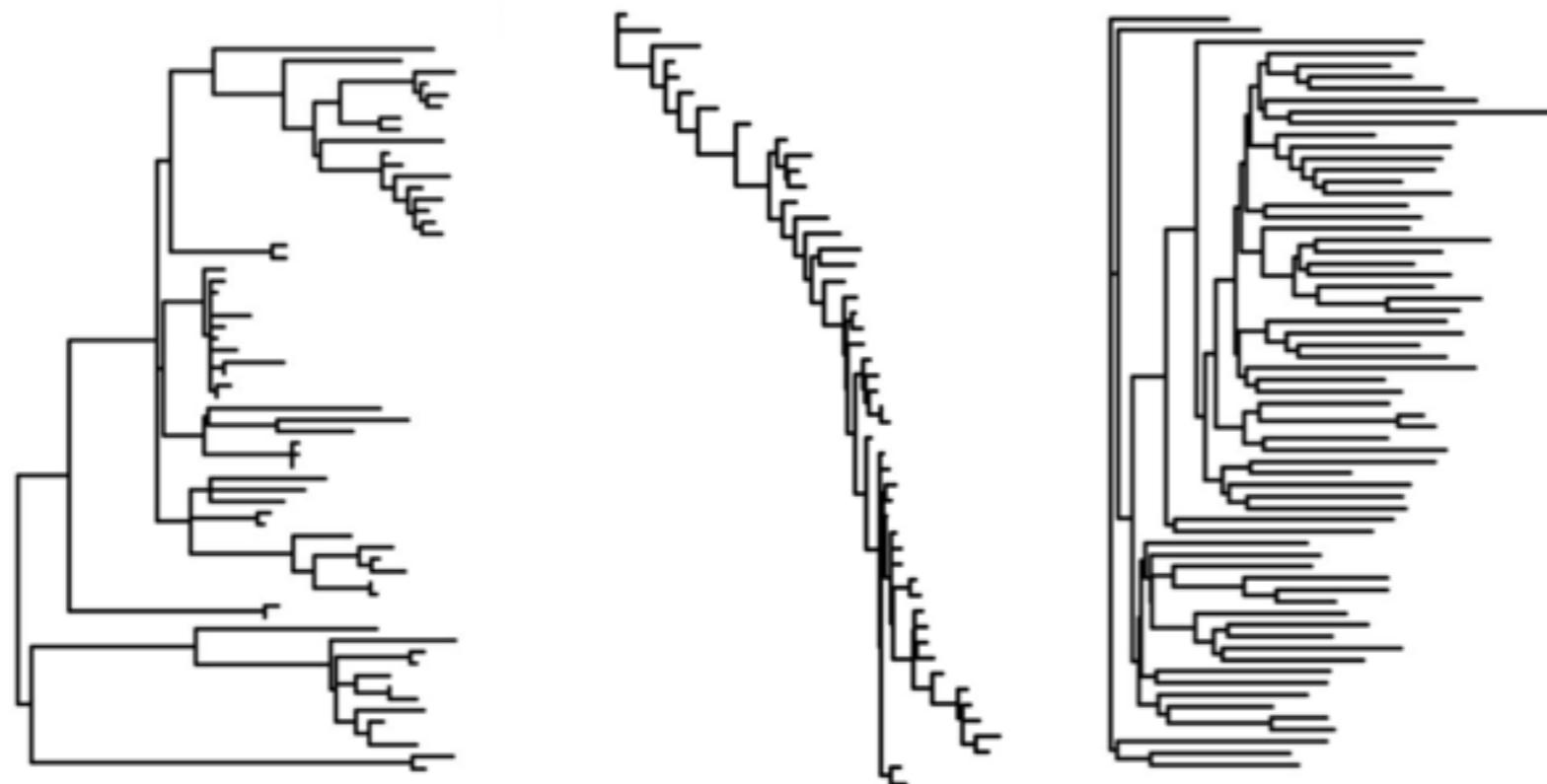
Swiss Institute of  
Bioinformatics

Veronika Bošková, Louis du Plessis, Joëlle Barido-Sottani  
Computational Evolution ([www.bsse.ch/cevo](http://www.bsse.ch/cevo))  
Department of Biosystems Science and Engineering



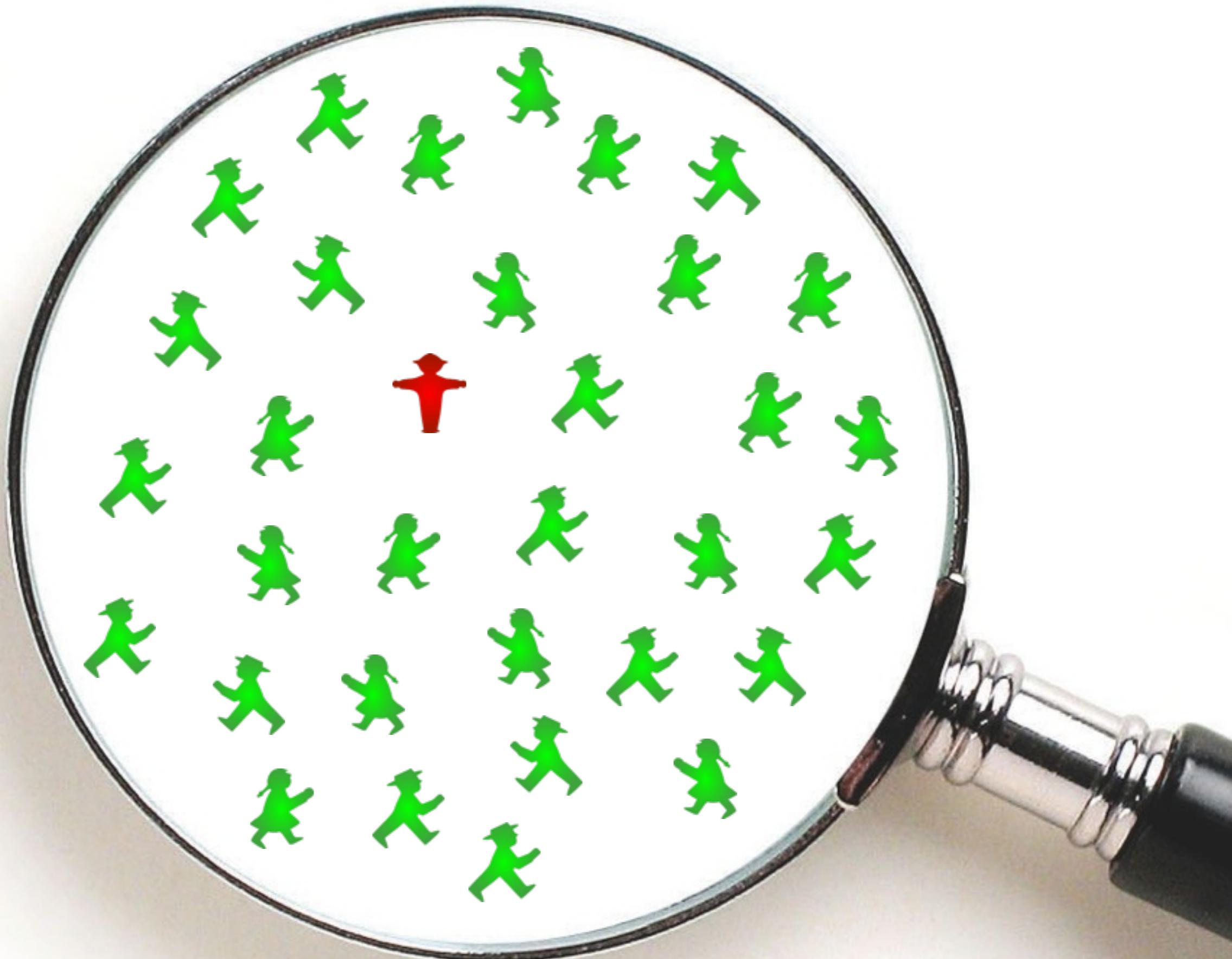


# Section 1: Introduction to Phylogenetics



---

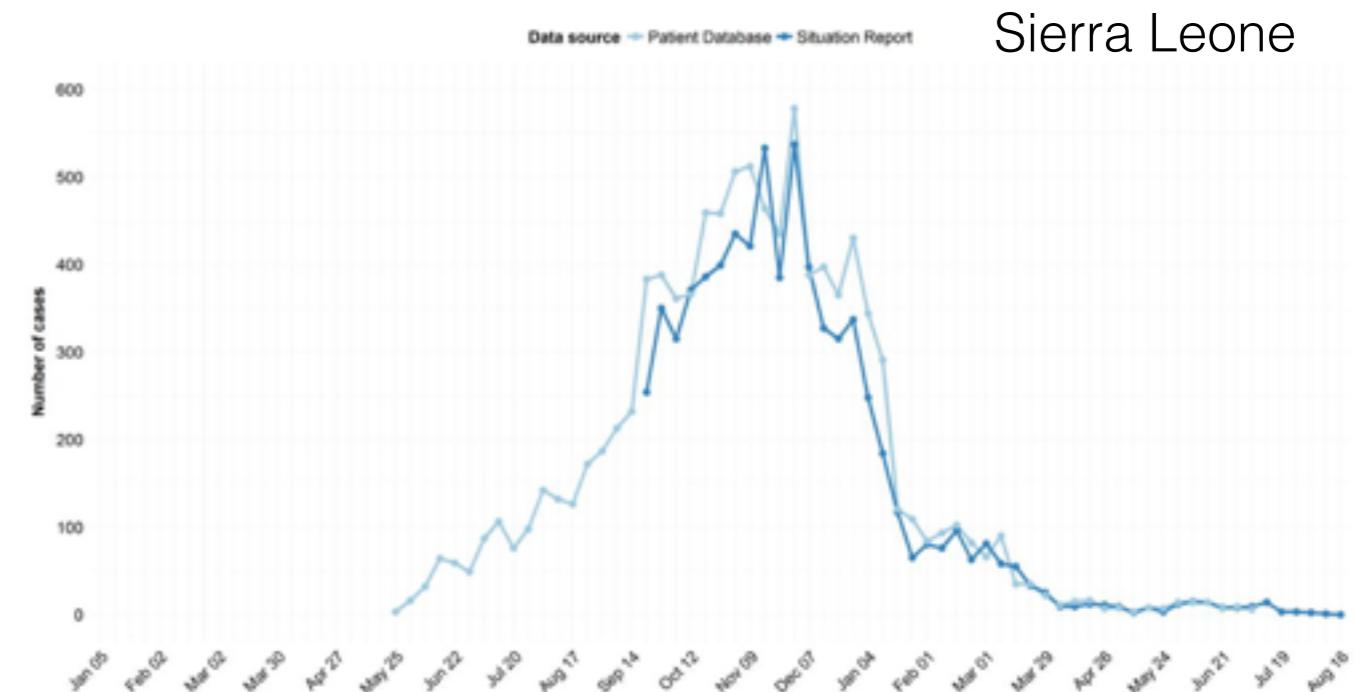
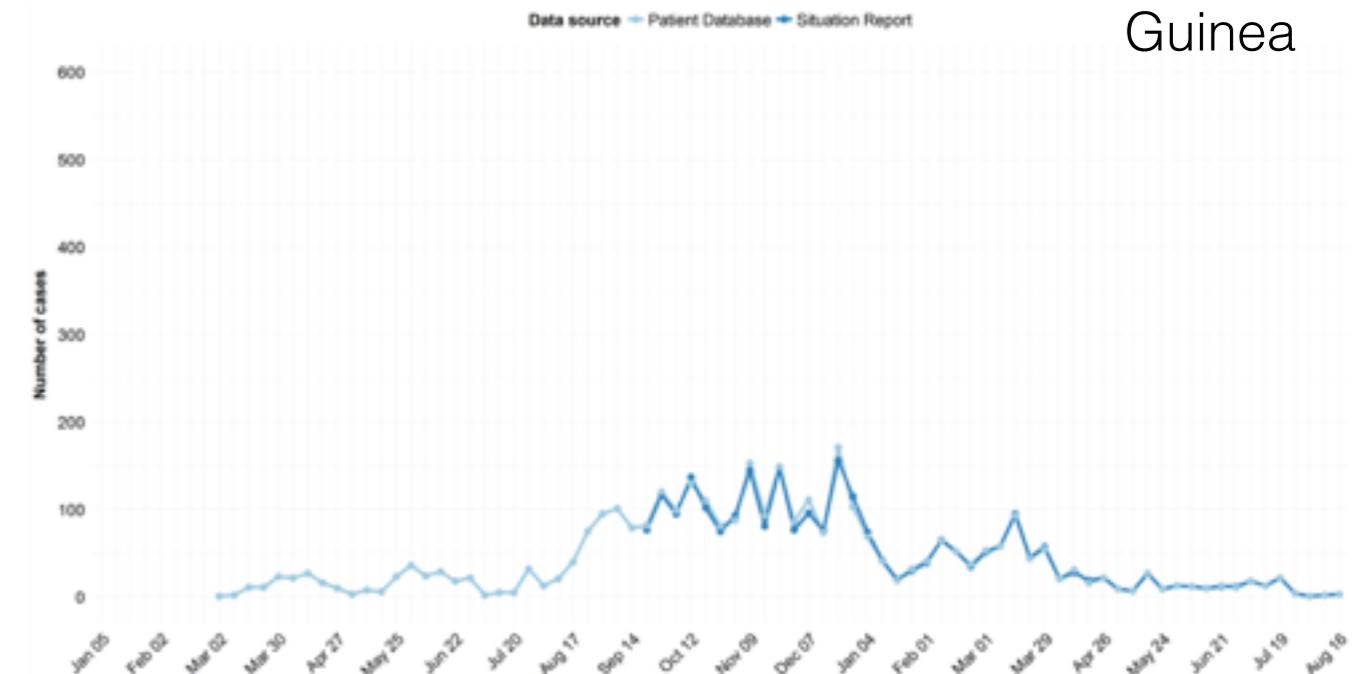
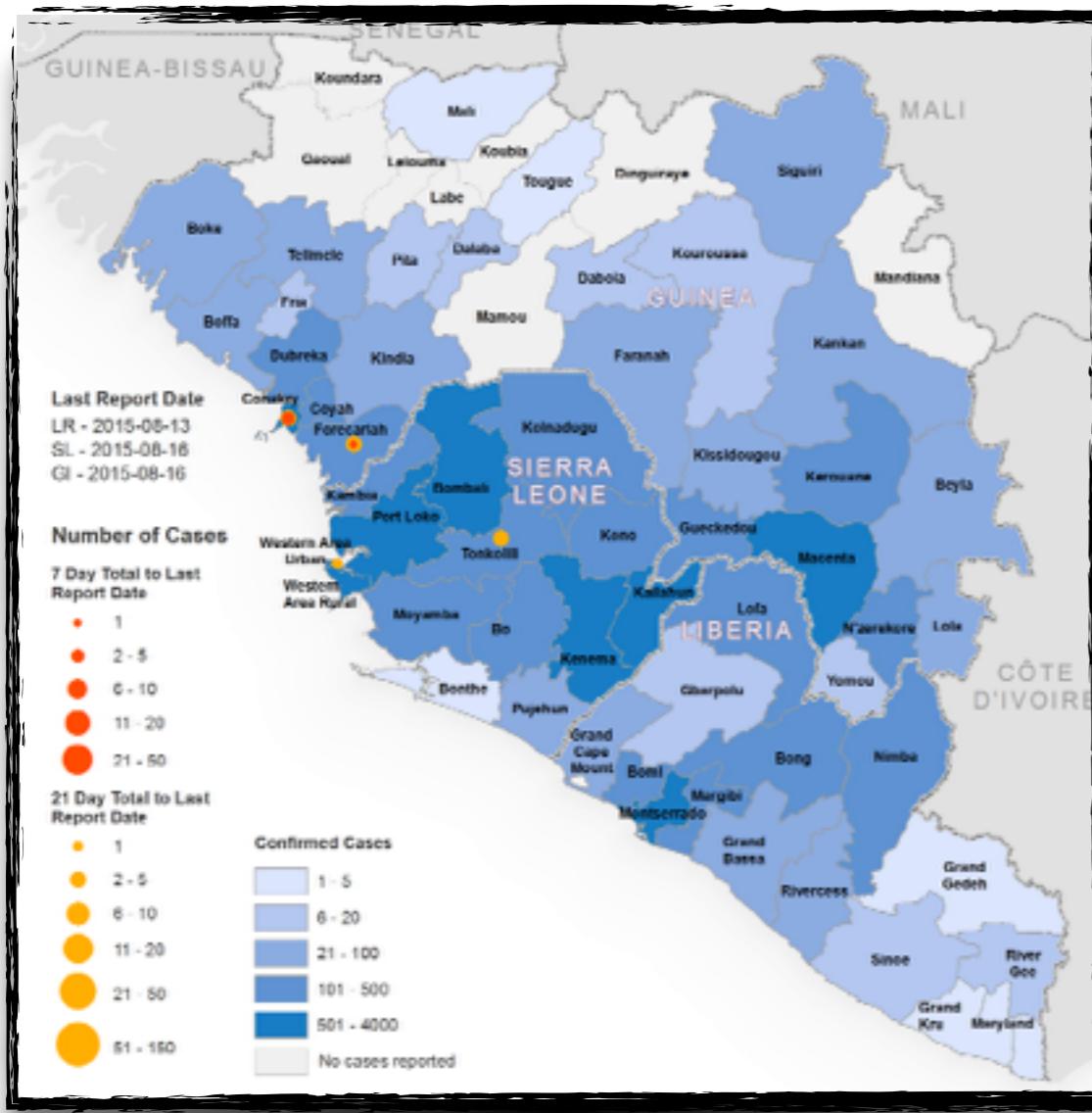
Louis du Plessis



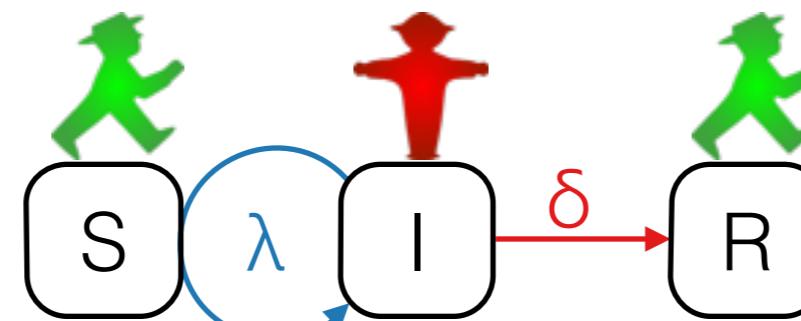
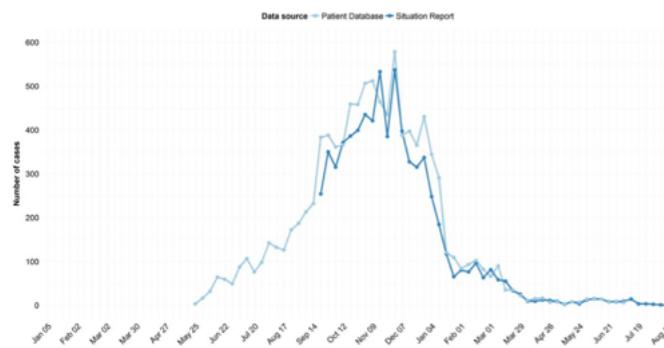




# How can we learn about epidemic spread?

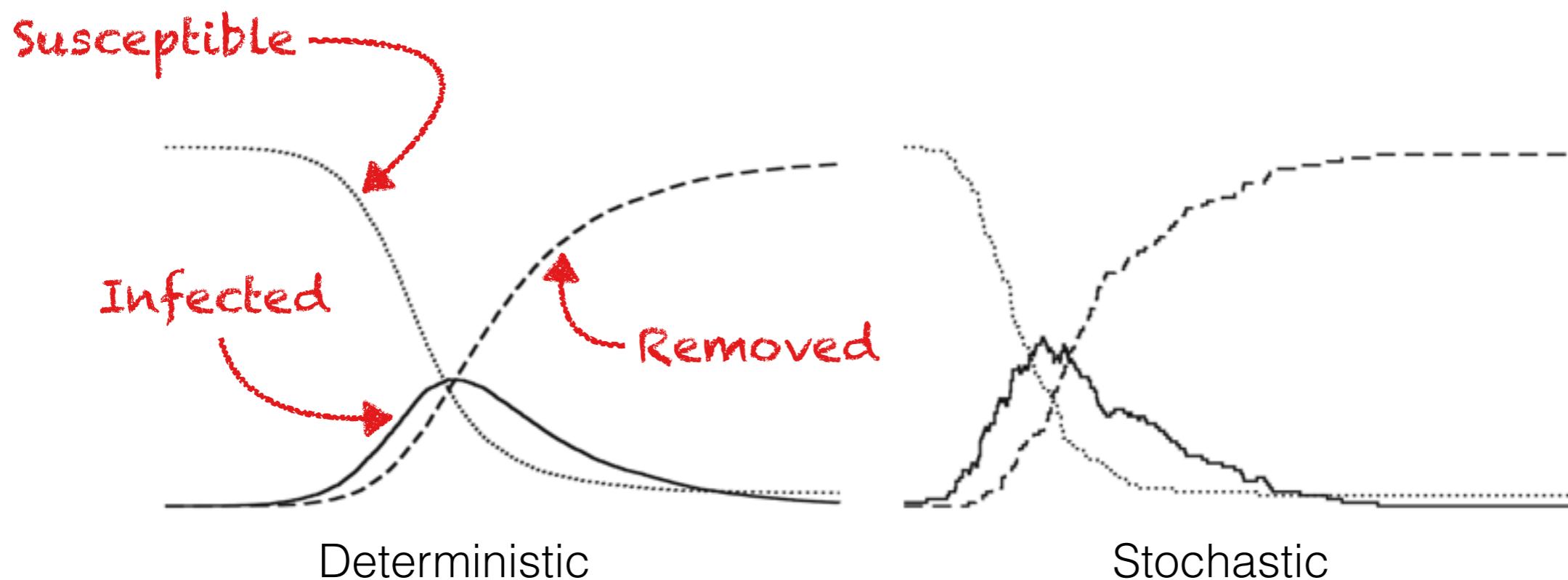


# Compartmental models for epidemics



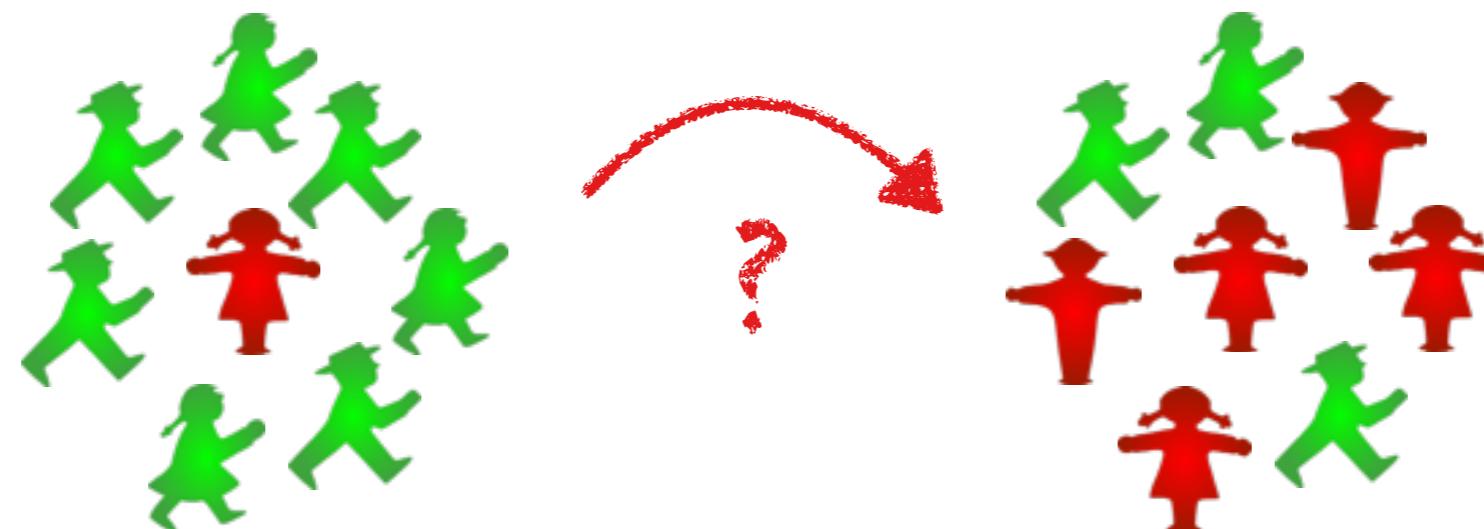
$$\begin{aligned} \frac{dS}{dt} &= -\lambda \frac{IS}{N} \\ \frac{dI}{dt} &= \lambda \frac{IS}{N} - \delta I \\ \frac{dR}{dt} &= \delta I \end{aligned}$$

- $\lambda$  — infection rate
- $\delta$  — becoming-noninfectious rate



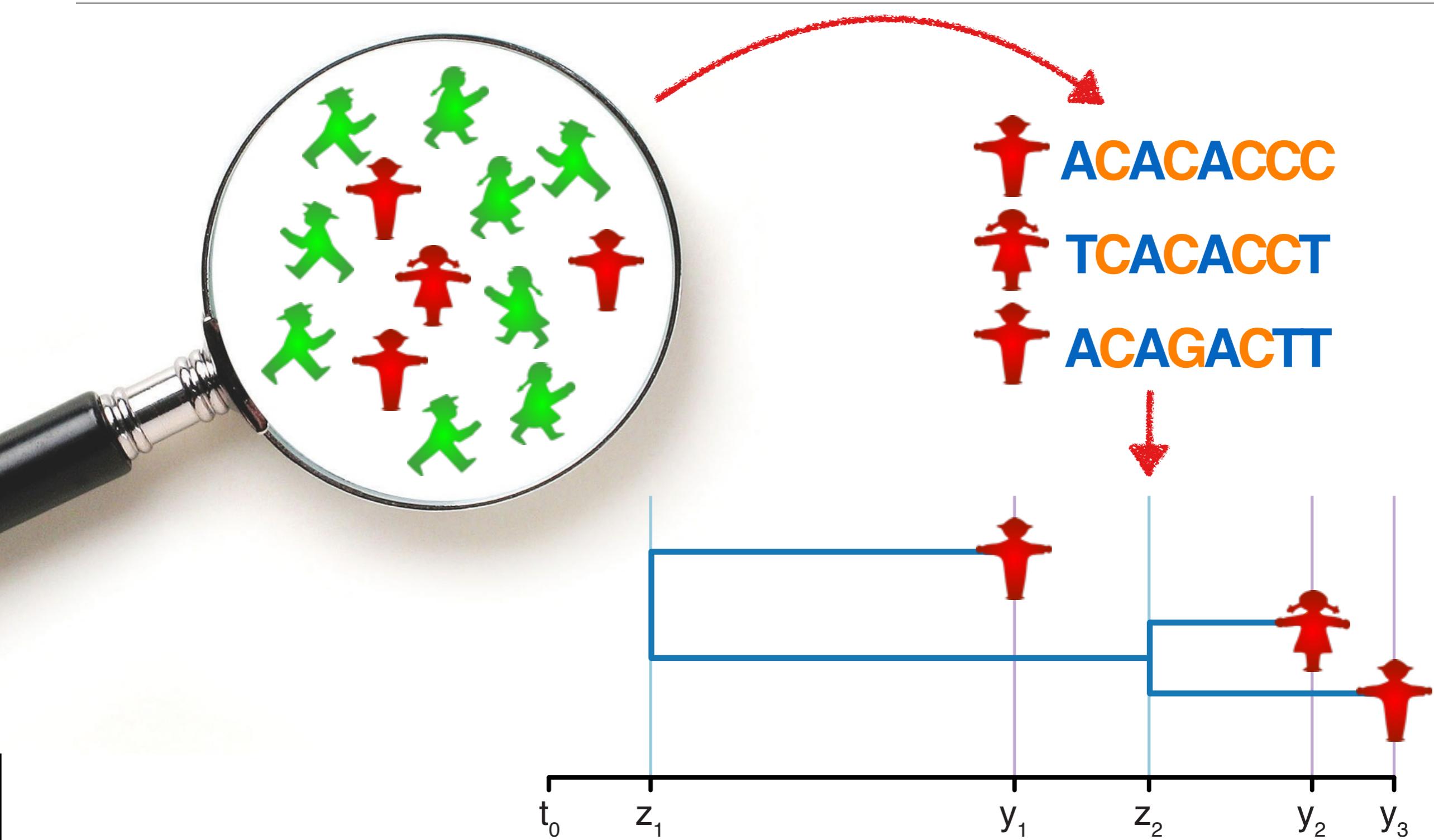
# What are the important questions?

---



- How fast is the epidemic spreading?  
(Basic reproduction number —  $R_0$ )
- When/where did the epidemic originate?
- Are public health interventions effective?  
(Effective reproduction number —  $R_e$ )
- Does population structure play a role?
- Did we sample everyone who is infected?
- ...?

# Can phylogenetics help?



# Can phylogenetics help?

## Transmission trees from sequencing data

Need a measurably evolving population

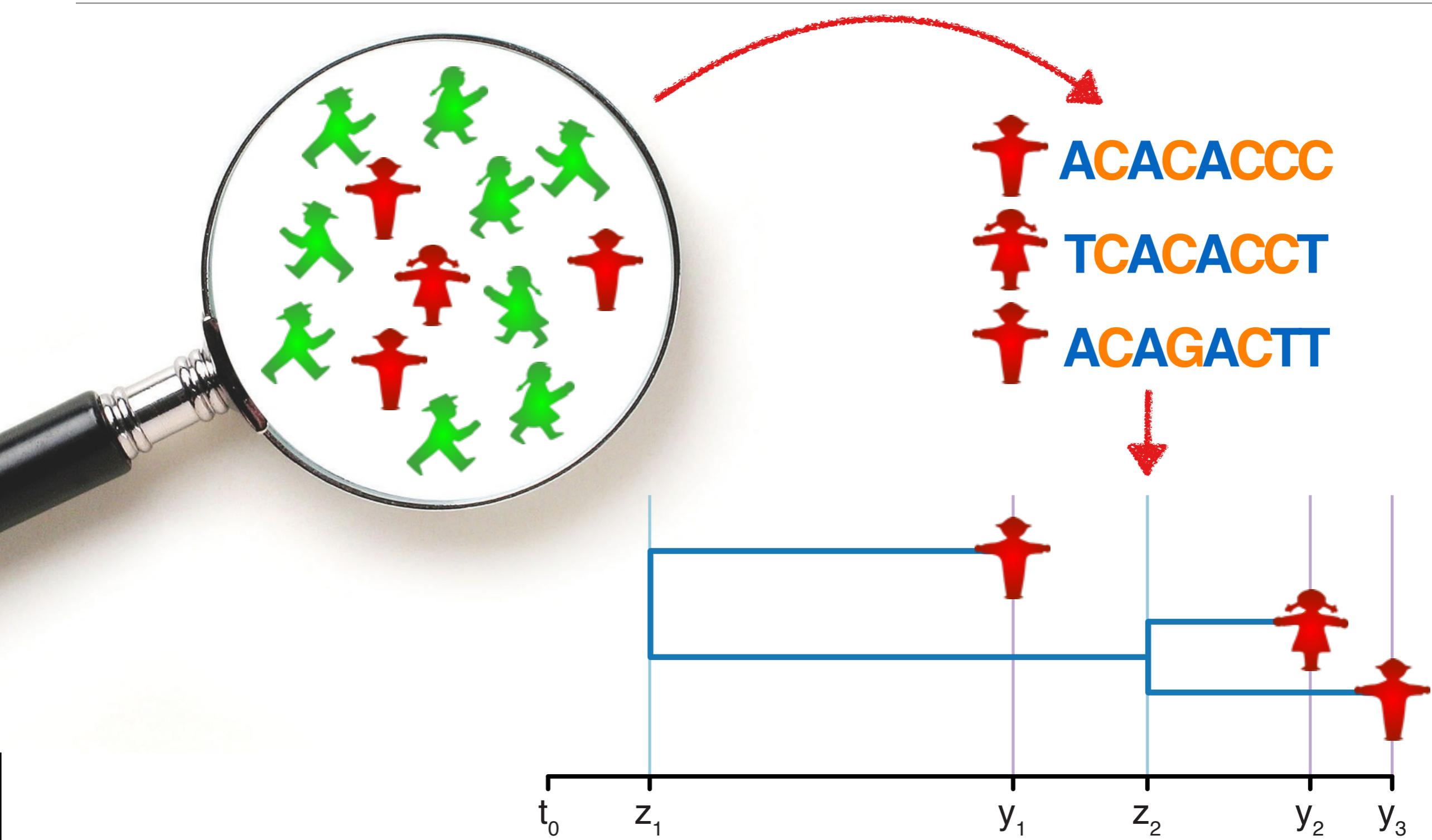
- Large population
- High mutation rate
- Short generation times

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

Ok... but is this useful?

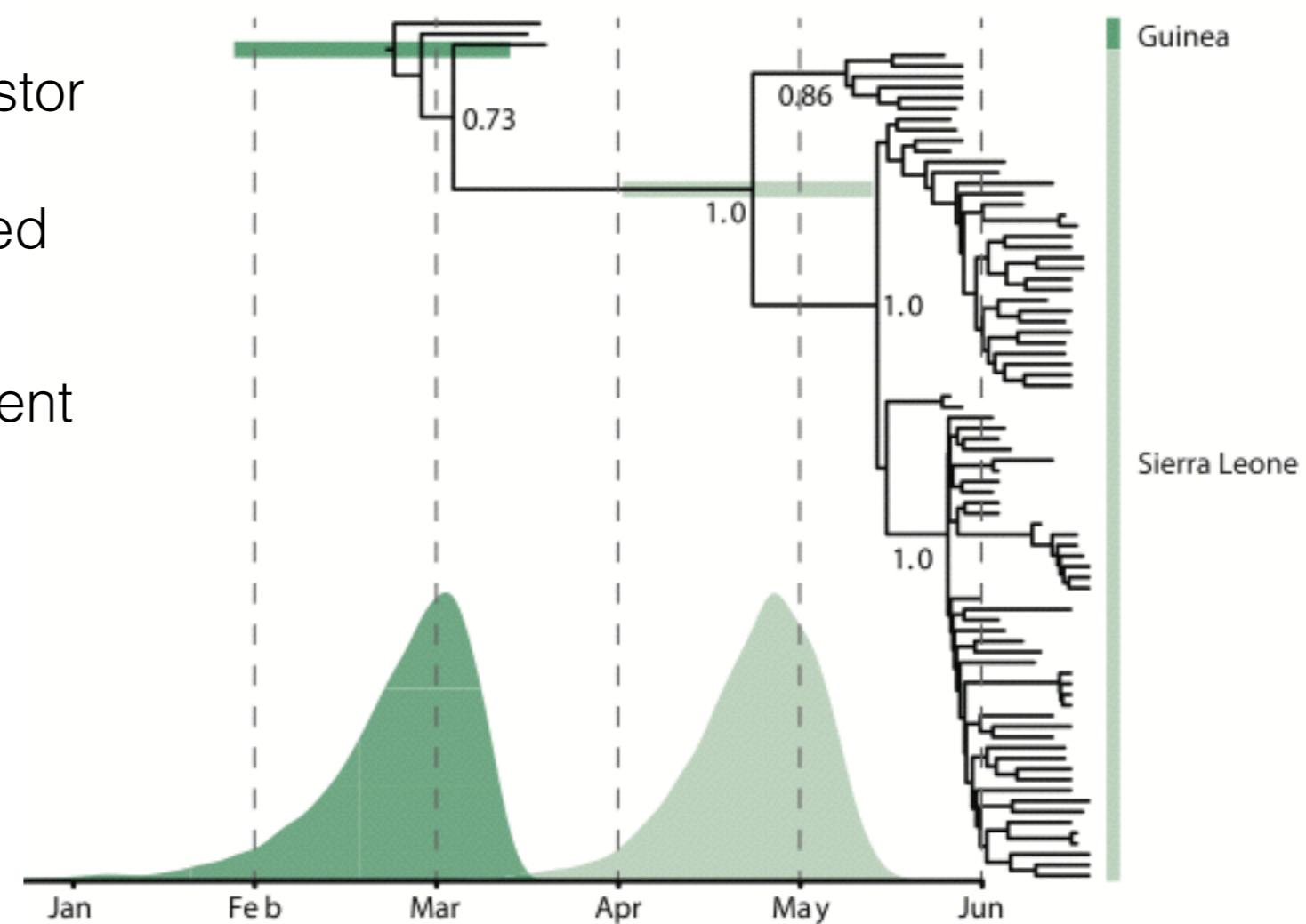
y<sub>3</sub>

# Can phylogenetics help?

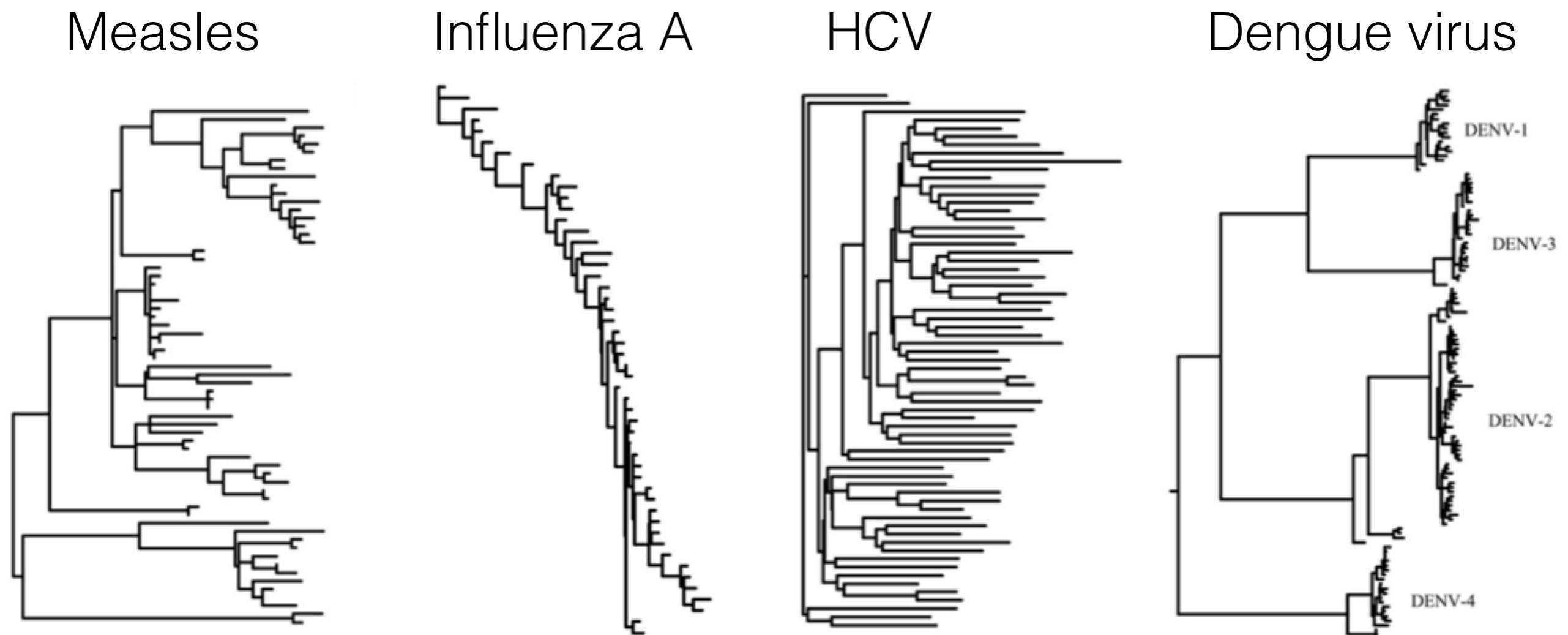


# Tracing epidemic origins: Ebola in West Africa

- Most recent common ancestor (MRCA) in February 2014 (2 months after first recorded case)
- Likely only one zoonotic event
- Similar techniques used to date the origin of the HIV pandemic and identify the zoonotic host

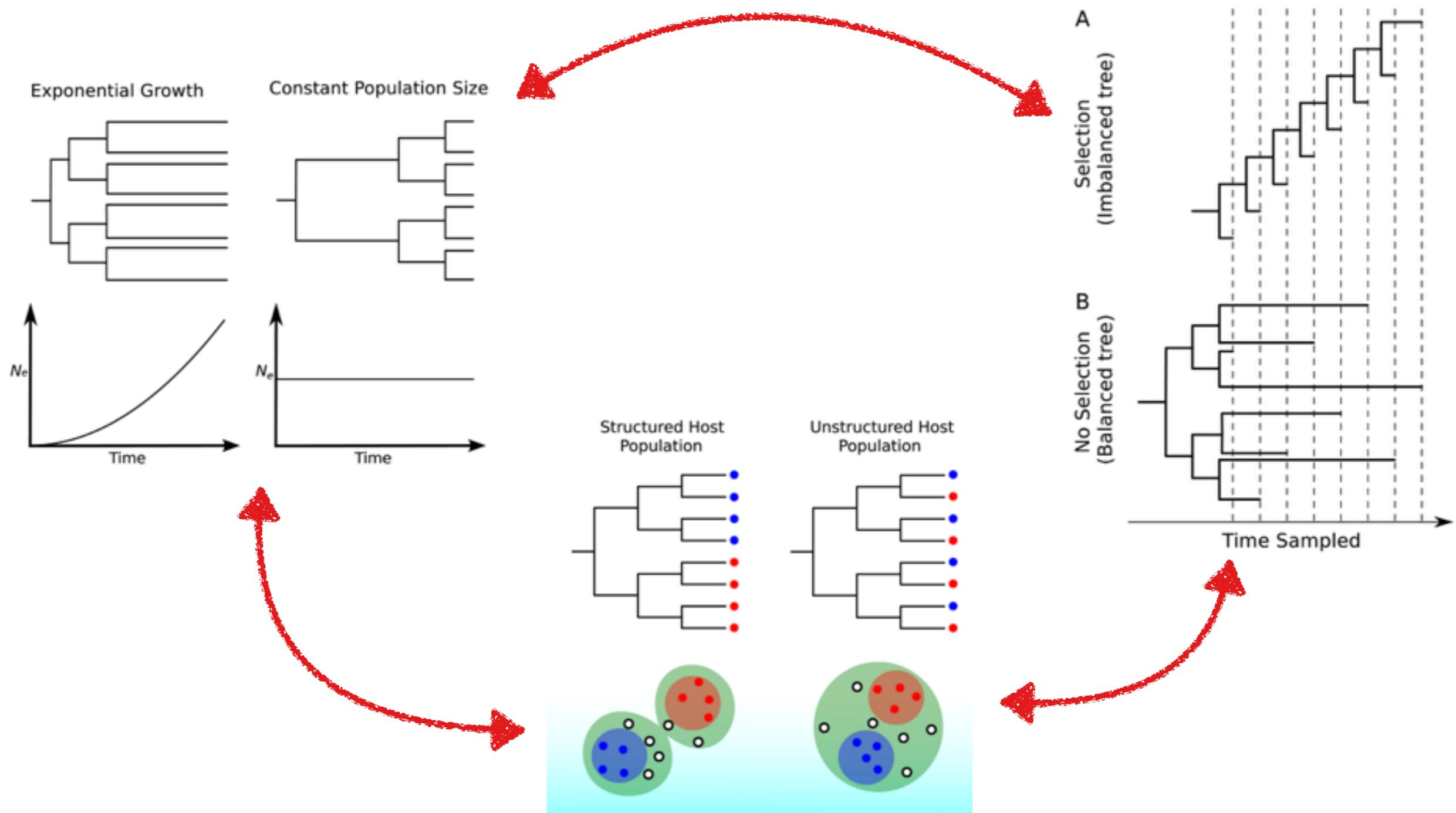


# Transmission trees from sequencing data



Why do the trees look  
so different?

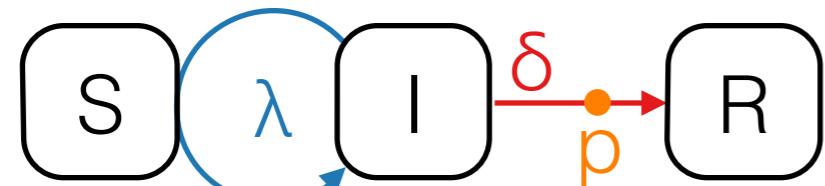
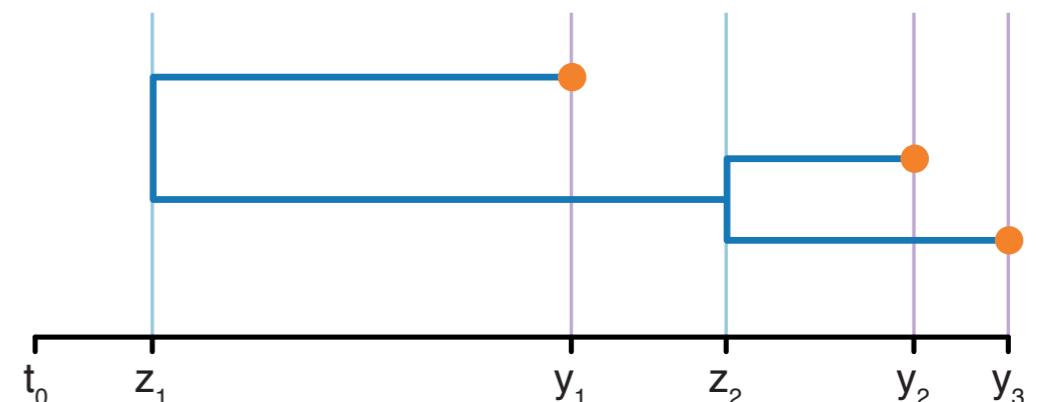
# Genomes contain a signature of the epidemiological dynamics



# What is phylodynamics?

## Phylogenetics

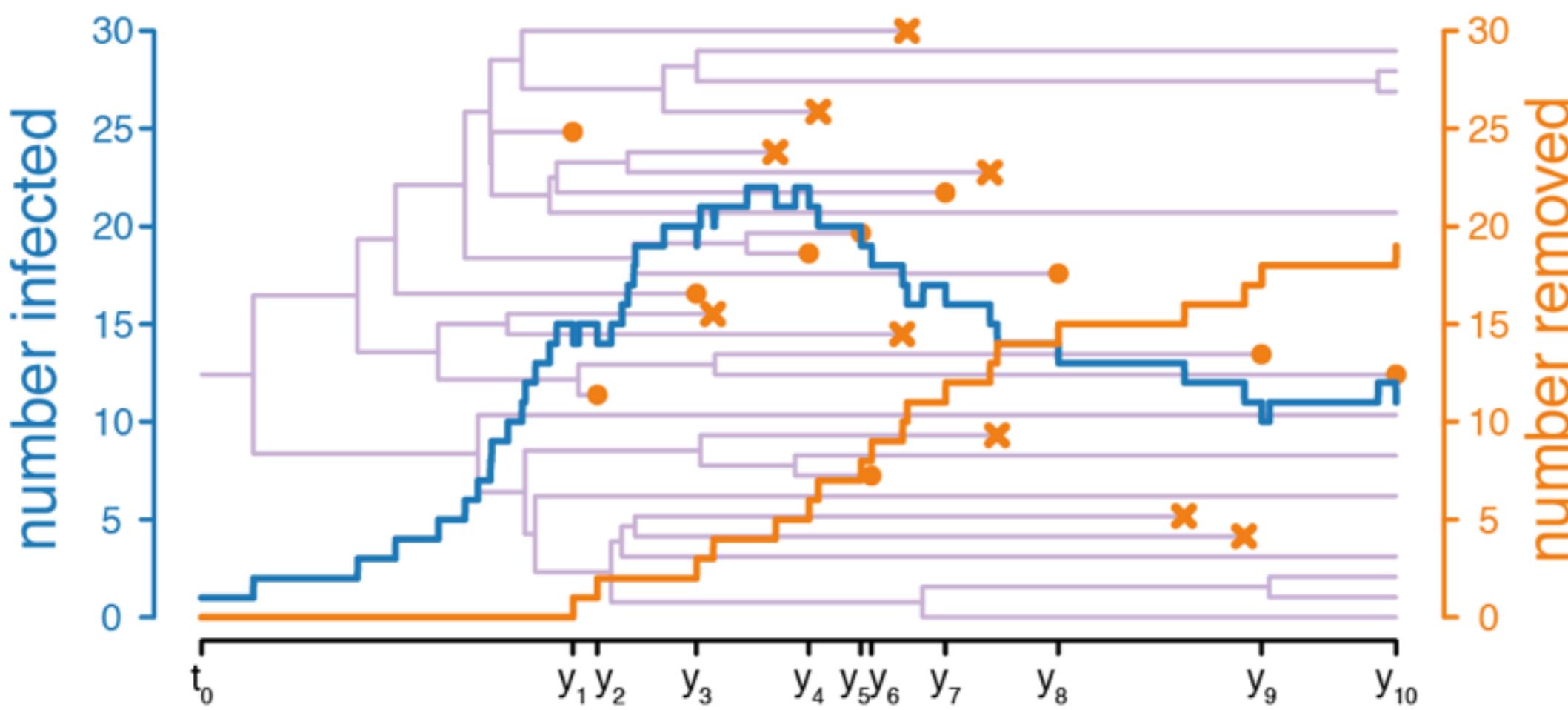
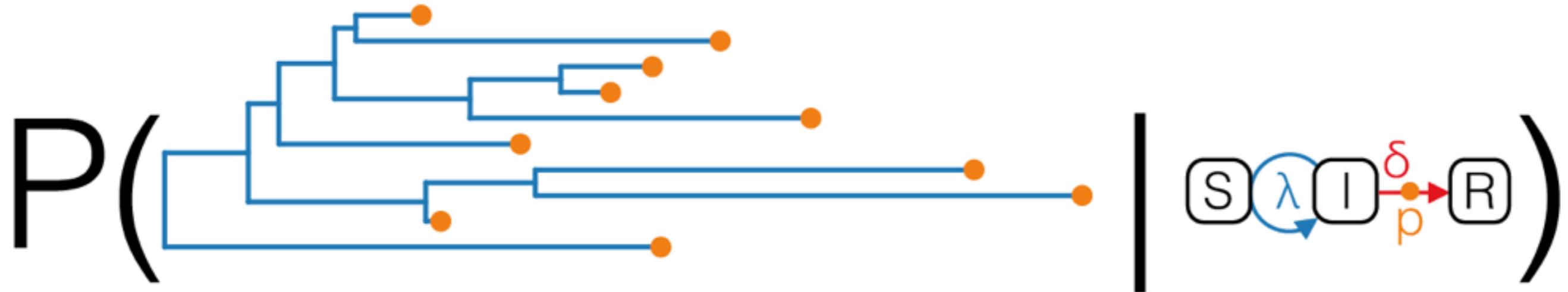
- State of process
- Transmission tree



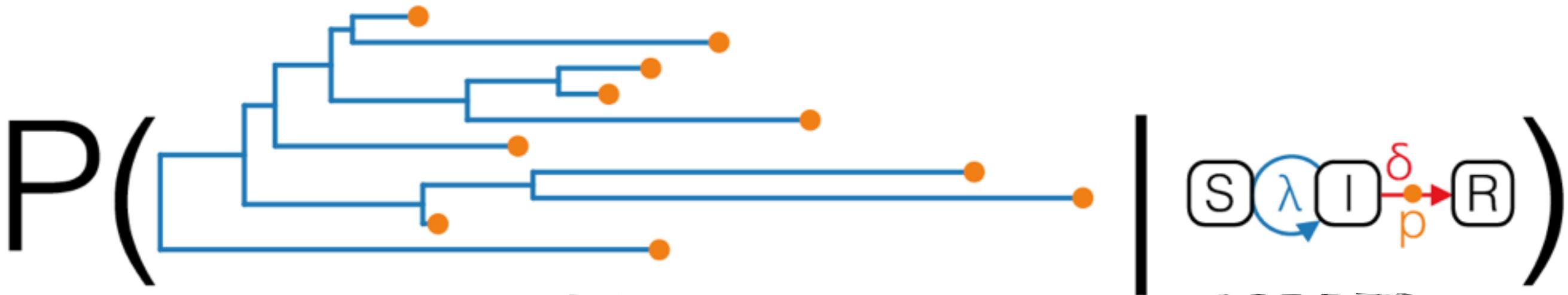
## Phylodynamics

- Dynamics of process
- Epidemiological parameters

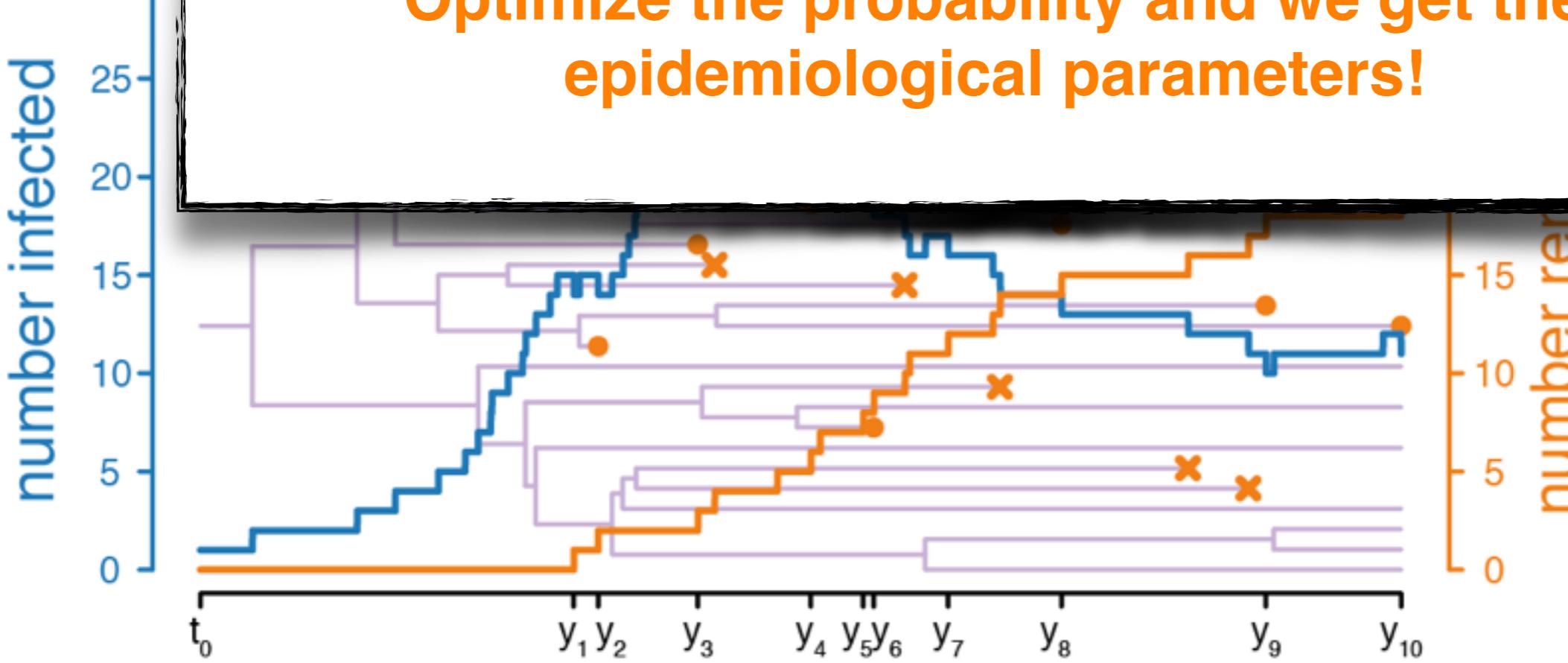
# A framework for phylodynamic inference



# A framework for phylodynamic inference



Optimize the probability and we get the epidemiological parameters!



# A framework for phylodynamic inference

---

## Step 1: Infer sampled transmission tree

Input data:

**ACACACCC**  
**TCACACCT**  
**TCAGACTT**  
**ACAGACTT**

 Sequence alignment

Model parameters:



Tree



Substitution model



Molecular clock model

Optimize:

$$P(\text{≡} \mid \text{Tree}, \text{Substitution model}, \text{Molecular clock model})$$

Find model parameters that best describe the evolution of sequences along a tree

## Step 2: Infer dynamics that led to the tree

Input data:



Tree  
(Realisation of a stochastic process)

Model parameters:



Demographic model

(Model describes dynamics of pathogen transmission)

Optimize:

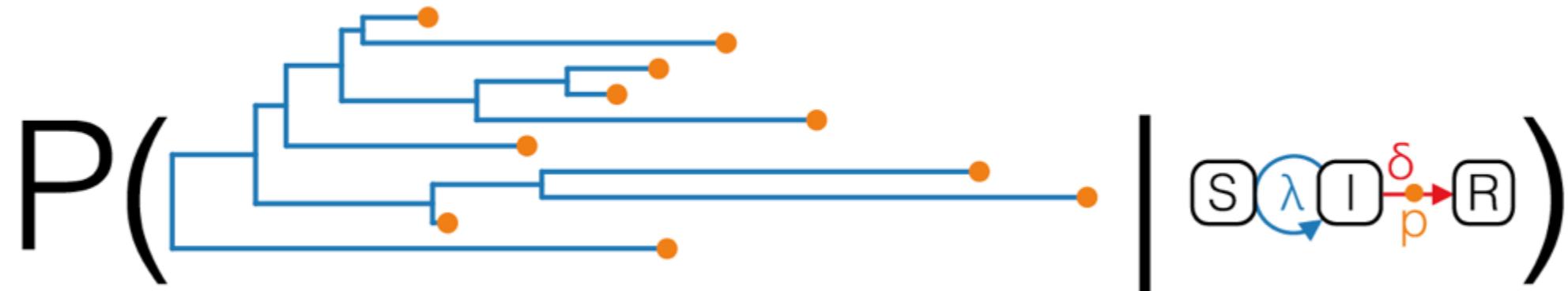
$$P(\text{Tree} \mid \text{Demographic model})$$

Find model parameters that best describe how the tree grows during an epidemic

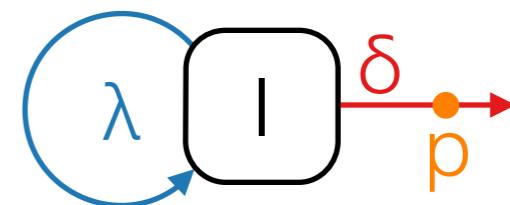
**Major assumption:** Patients do not transmit after they are sampled!

# Birth-death models (forward in time)

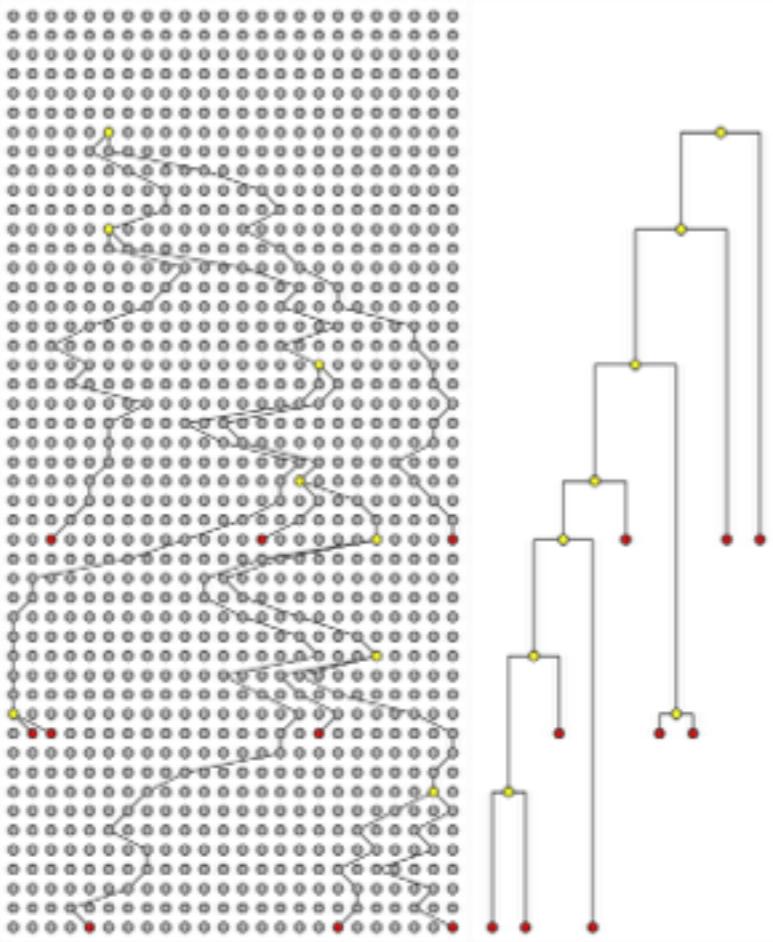
---



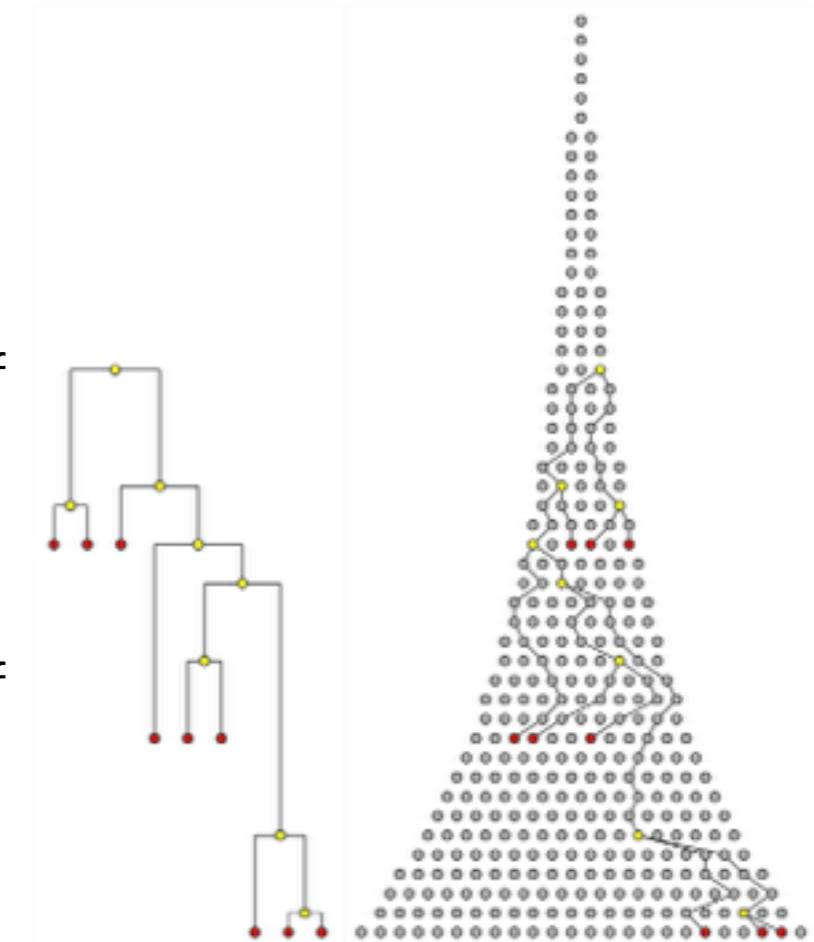
- Directly calculate the probability of observing a sampled tree given a model
- Computationally expensive for more complex models
- Cannot calculate probability in general
- Usually stick with simple models:



# Coalescent approximation (backward in time)



- Relates the shape of the tree to population size
  - Assume a random sample drawn from a population
  - Calculate the probability of two lineages coalescing in a given time under a Wright-Fisher process
  - Calculate the probability of observing the coalescent times in the sampled tree
  - Infer effective population size ( $N_e$ ) or growth rate
- 
- Originally assumed a constant population size
  - Generalized to any deterministically changing population size  
(eg. exponential growth, SIR dynamics)



# Summary: Phylodynamic inference framework

---

Epidemiological and evolutionary dynamics on the same timescale



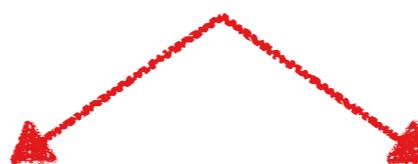
Pathogen genomes contain signature of epidemic spread



Leaves pattern on resulting phylogeny (sampled transmission chain)



Extract epidemiological parameters from phylogeny



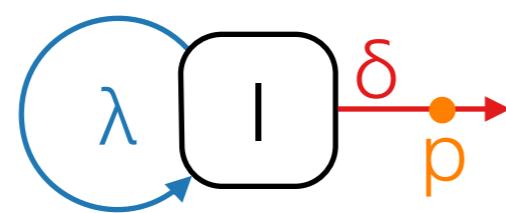
Birth-death models

Coalescent models

# Summary: Phylodynamic inference framework

## Epidemiological parameters example

Constant rate birth-death process



- $\lambda$  — infection rate
- $\delta$  — becoming-noninfectious rate
- $p$  — sampling probability

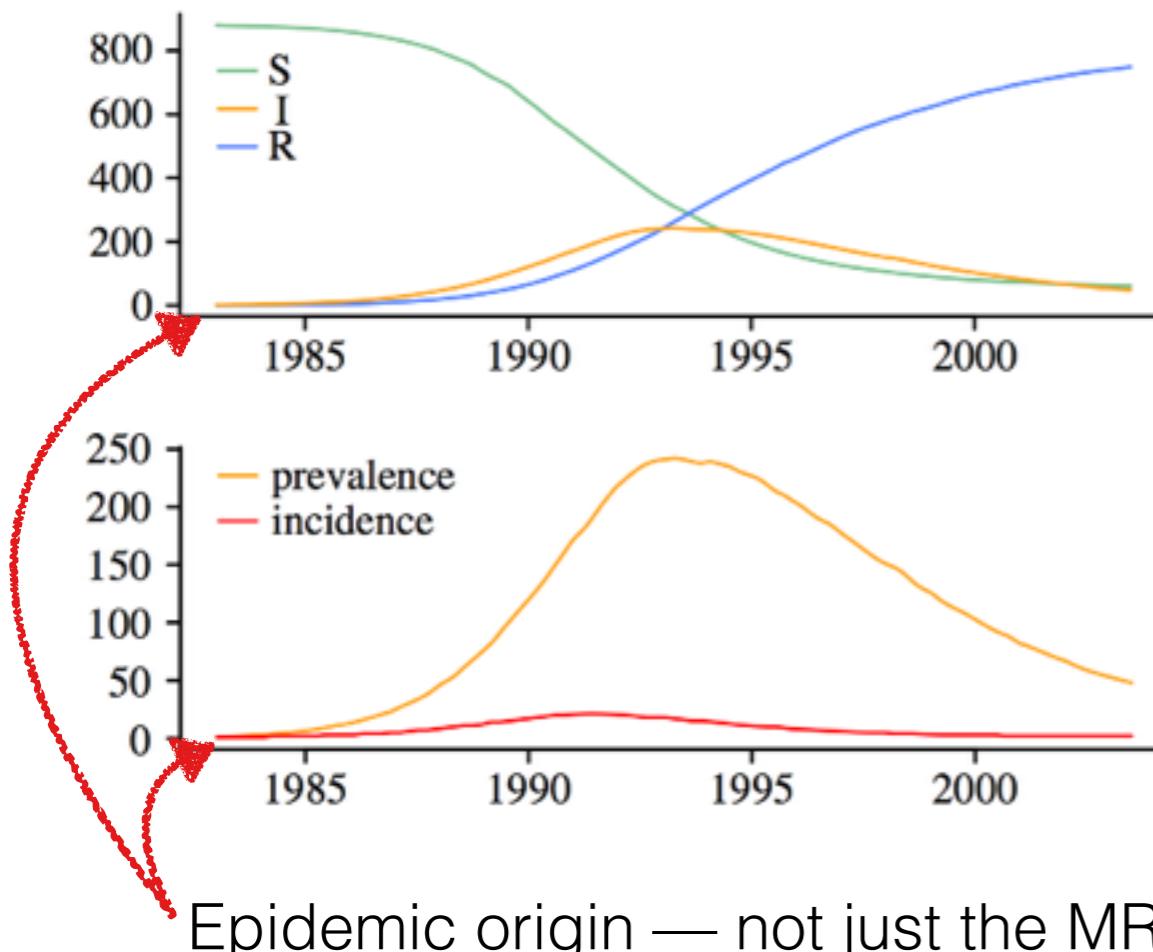
- $R_0 = R_e = \lambda/\delta$
- Infectious period =  $1/\delta$
- Epidemic origin: Time when  $I = 0$
- Number infected: Growing exponentially at rate  $\lambda - \delta$
- ...

Birth-death models

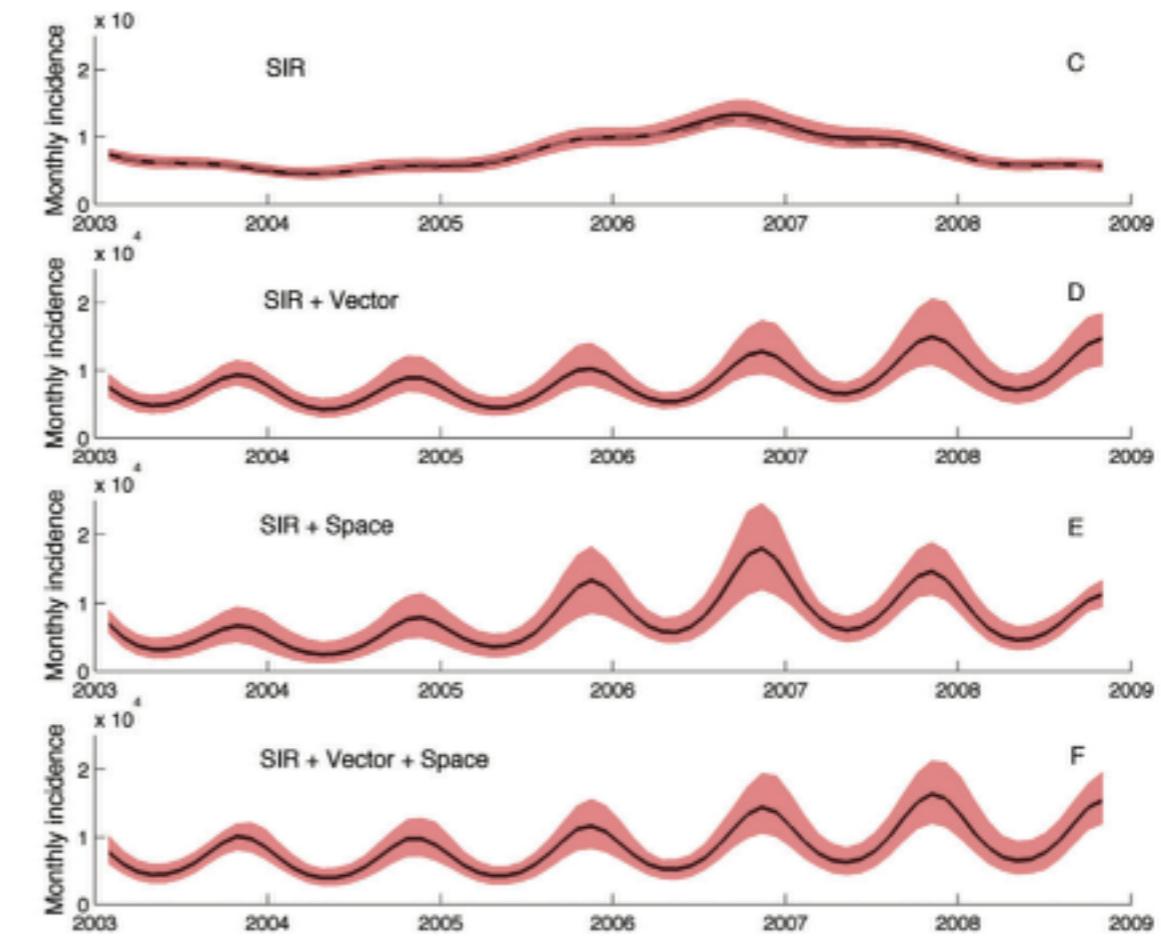
Coalescent models

# Mechanistic epidemiological models

Birth-death — HIV-1



Coalescent — DENV-1



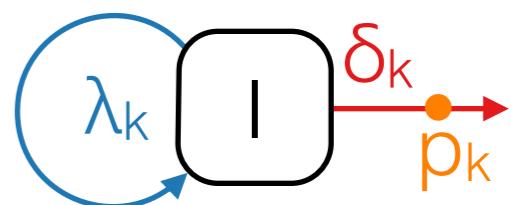
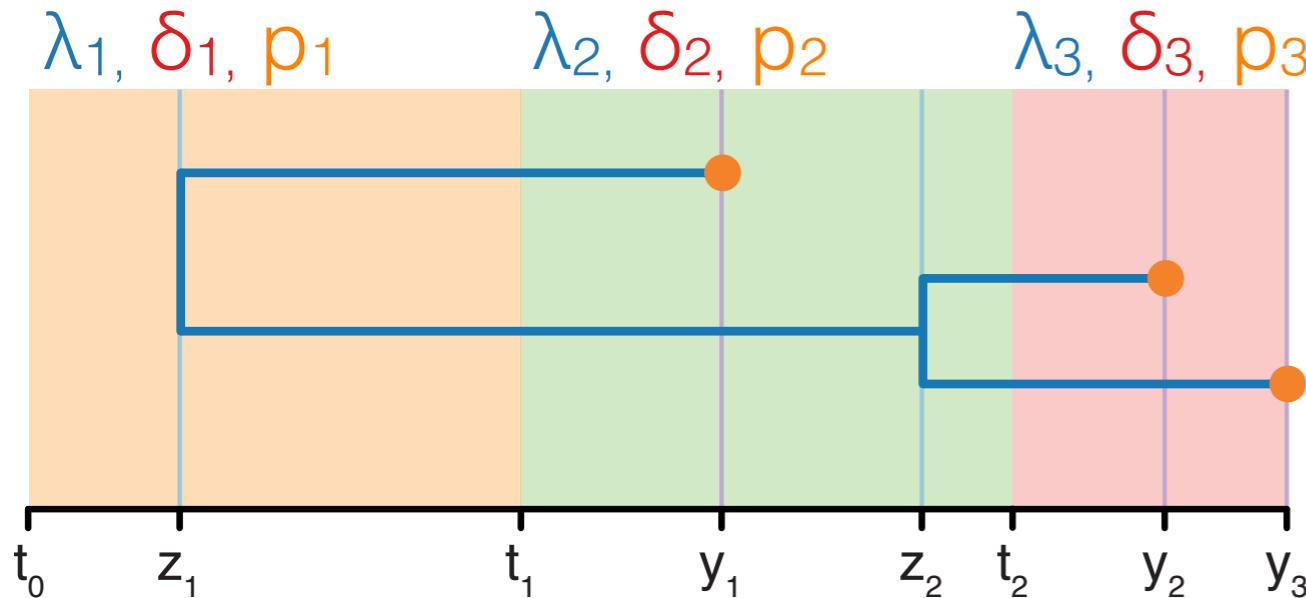
Epidemic origin — not just the MRCA!

**Easier to handle more complex mechanistic dynamics in a coalescent framework!**

# Nonparametric models: Skyline approximation

- Mechanistic models attempt to exactly describe the dynamics of the epidemic
- Not always possible to derive an expression for the likelihood under complex dynamics
- Exact dynamics are not always known

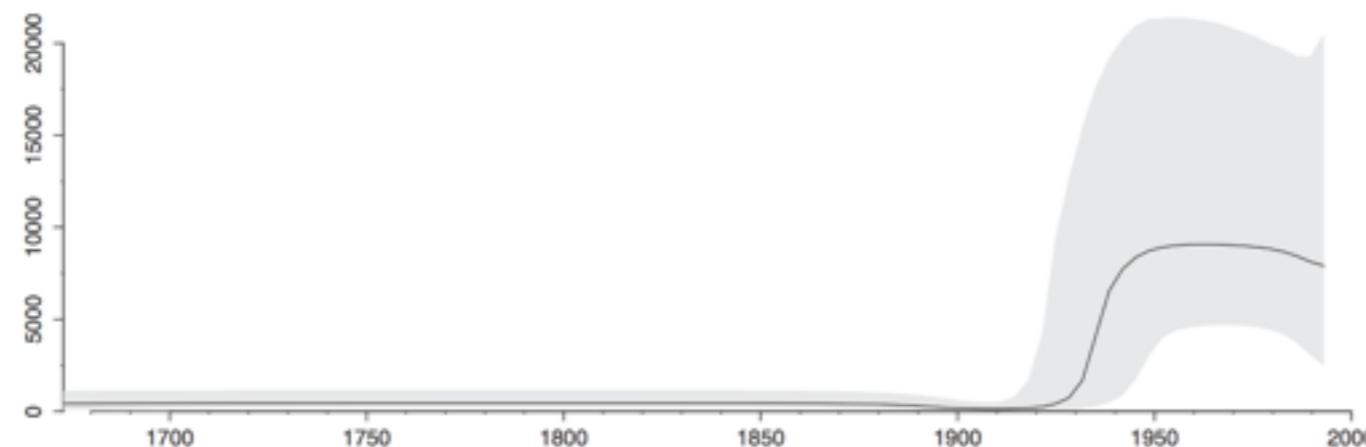
## Solution: Piecewise constant approximation



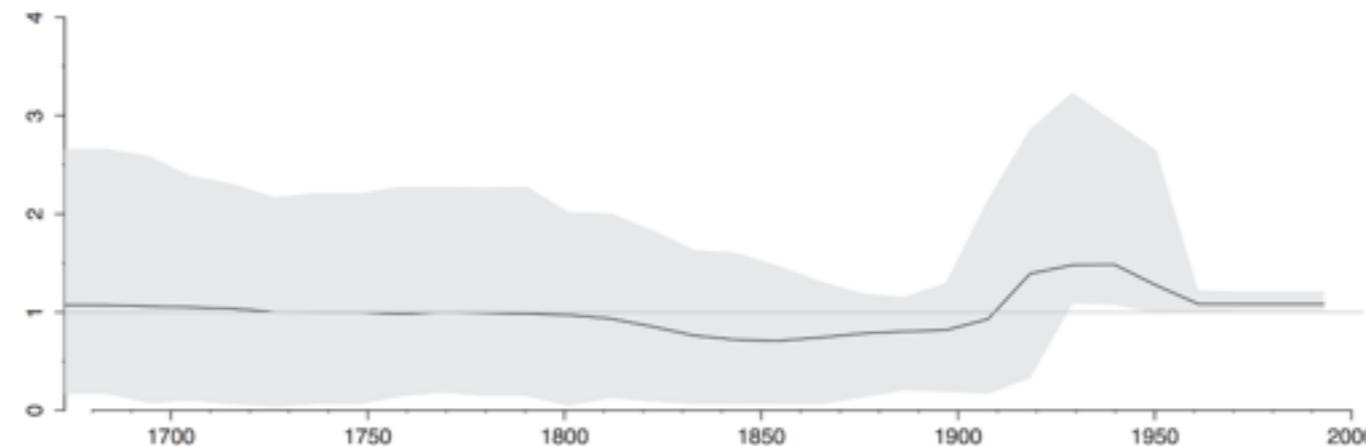
- Birth-death skyline:  
Rates change over time
- Coalescent skyline:  
 $N_e$  changes over time
- Can approximate any type of dynamics  
(if population structure not important)
- Need more data  
(no prior information on dynamics)

# Skyline approximation: HCV in Egypt

Coalescent Skyline — Changes in effective population size,  $N_e$



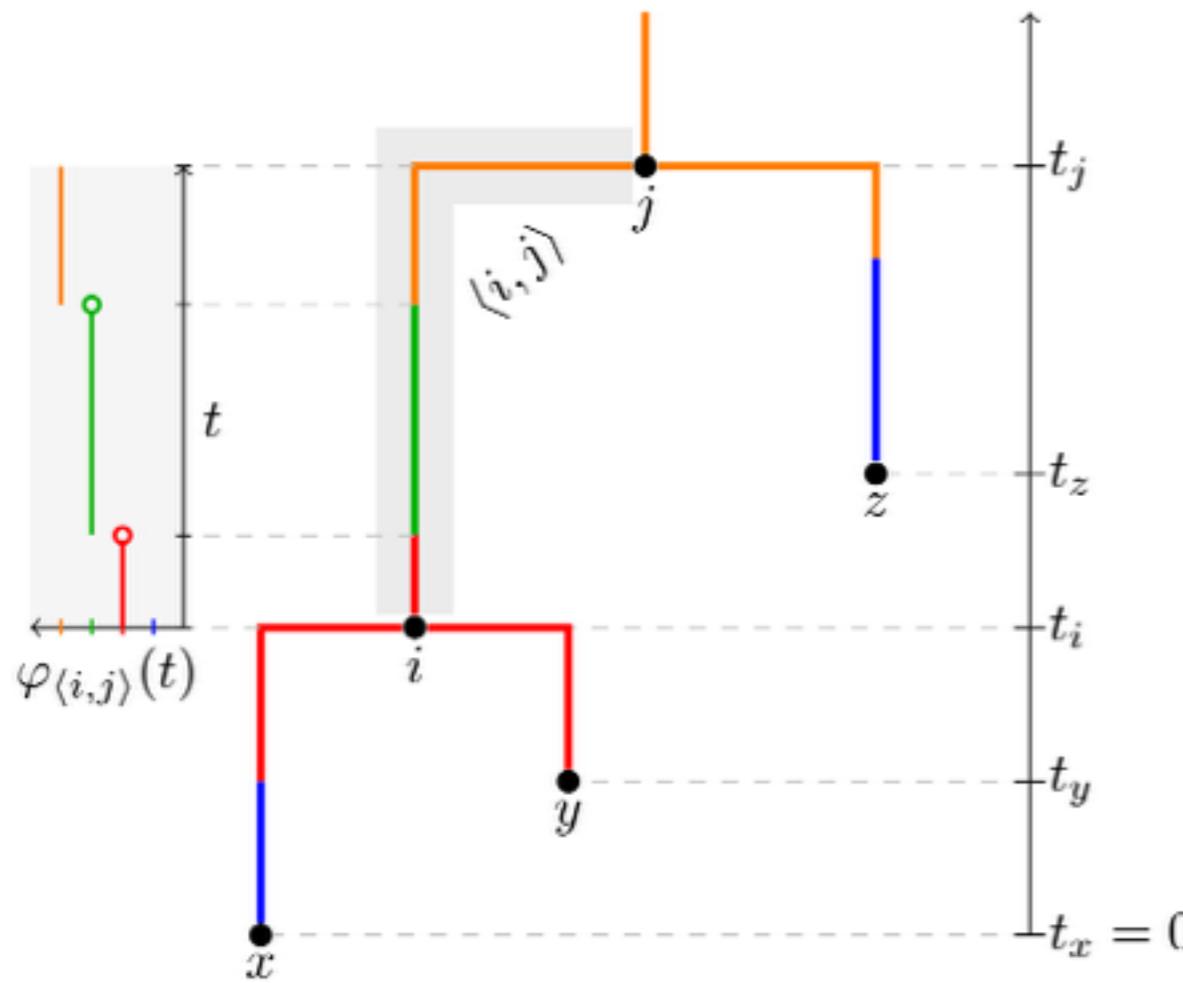
Birth-death Skyline — Changes in effective reproduction number,  $R_e$



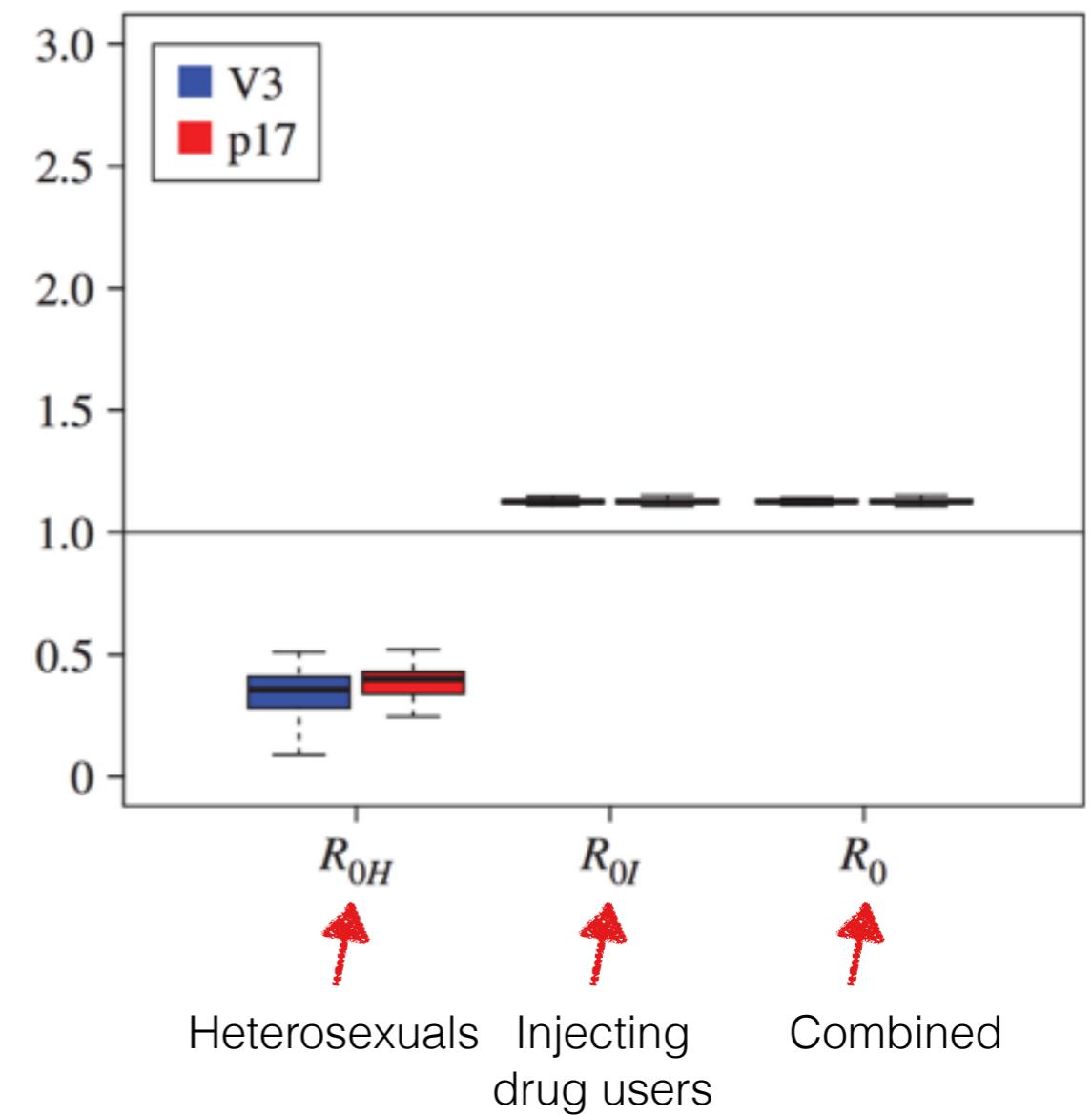
Large increase in  $N_e/R_e$  around 1950s due to antischistosomal injections with contaminated needles!

# Population structure

## HIV-1 in Latvia

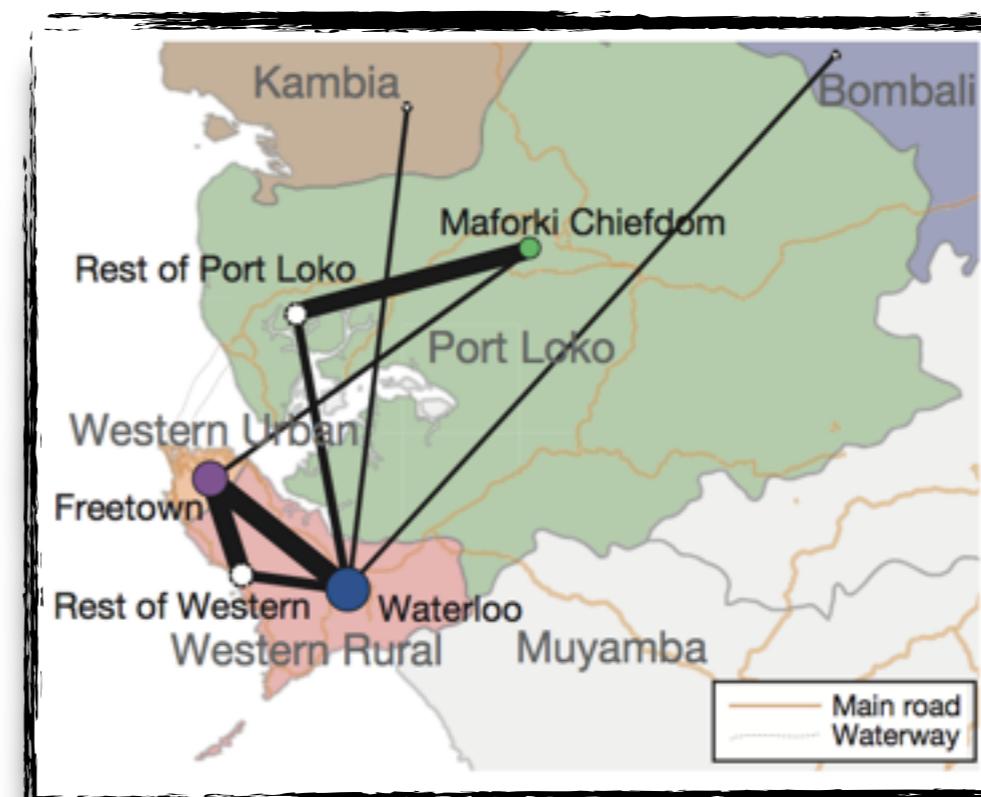
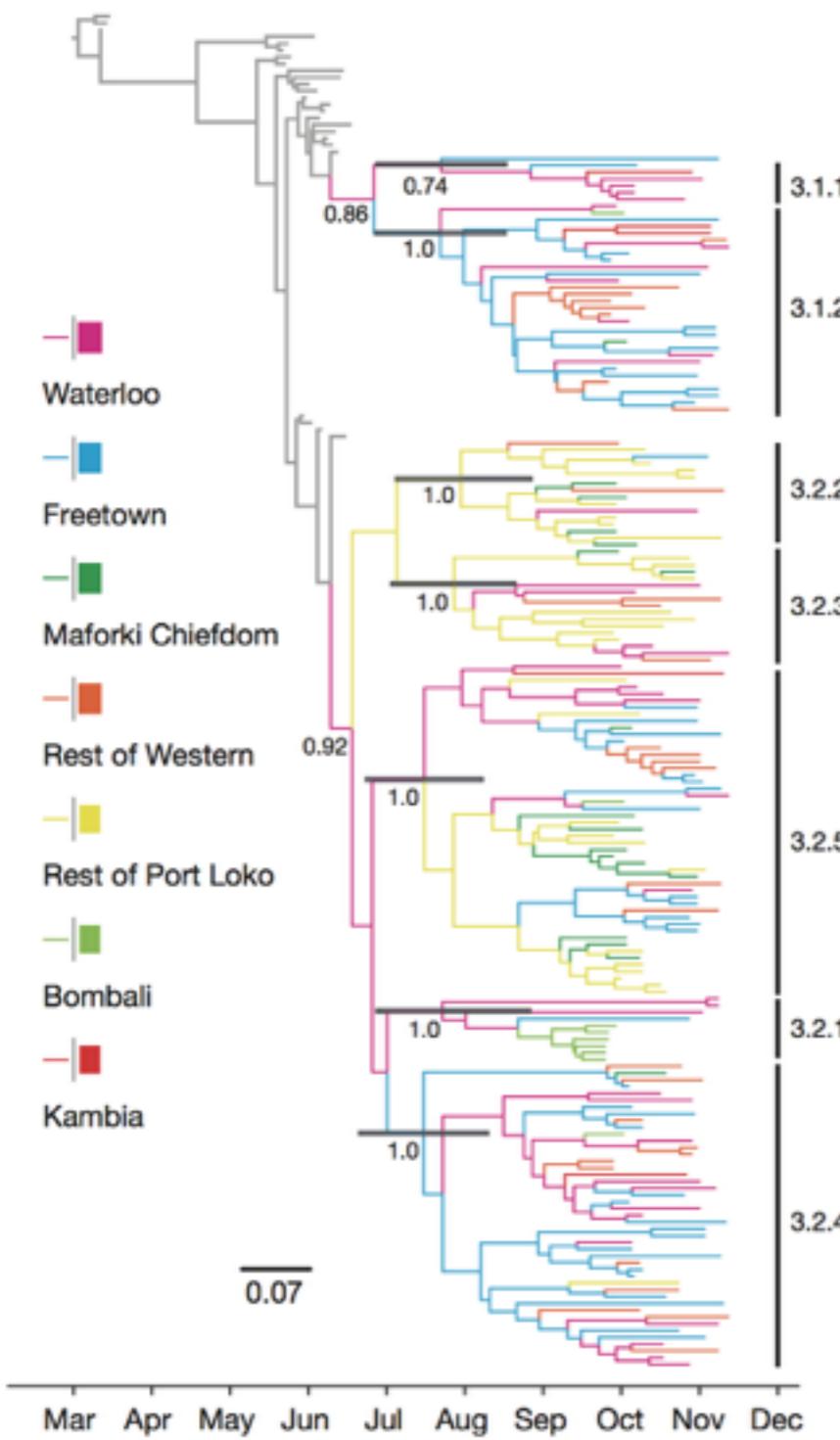


- Explicitly model different types of subpopulations
- Different types may be associated with different epidemiological parameters
- Can also model migration between types



- Injecting drug users have  $R_0 > 1$
- This is keeping the epidemic alive!

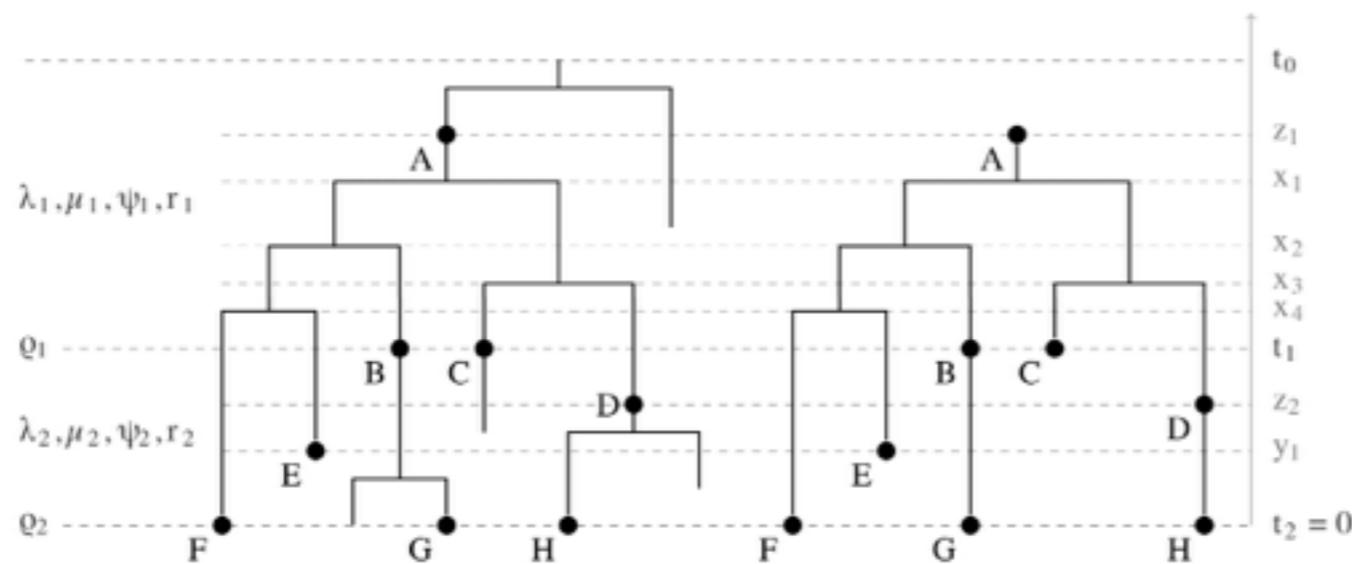
# Phylogeography: Ebola in West Africa



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

# Sampled ancestors

- Sometimes individuals keep on transmitting after they are sampled (eg. HIV-1, Ebola virus)
- Can relax this assumption for birth-death models



# Samp

## HIV-1

- Sometimes after the (eg. HIV)
- Can re models

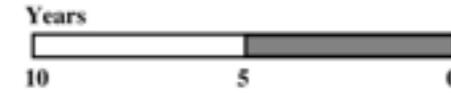
$\lambda_1, \mu_1, \psi_1, r_1$

$Q_1$

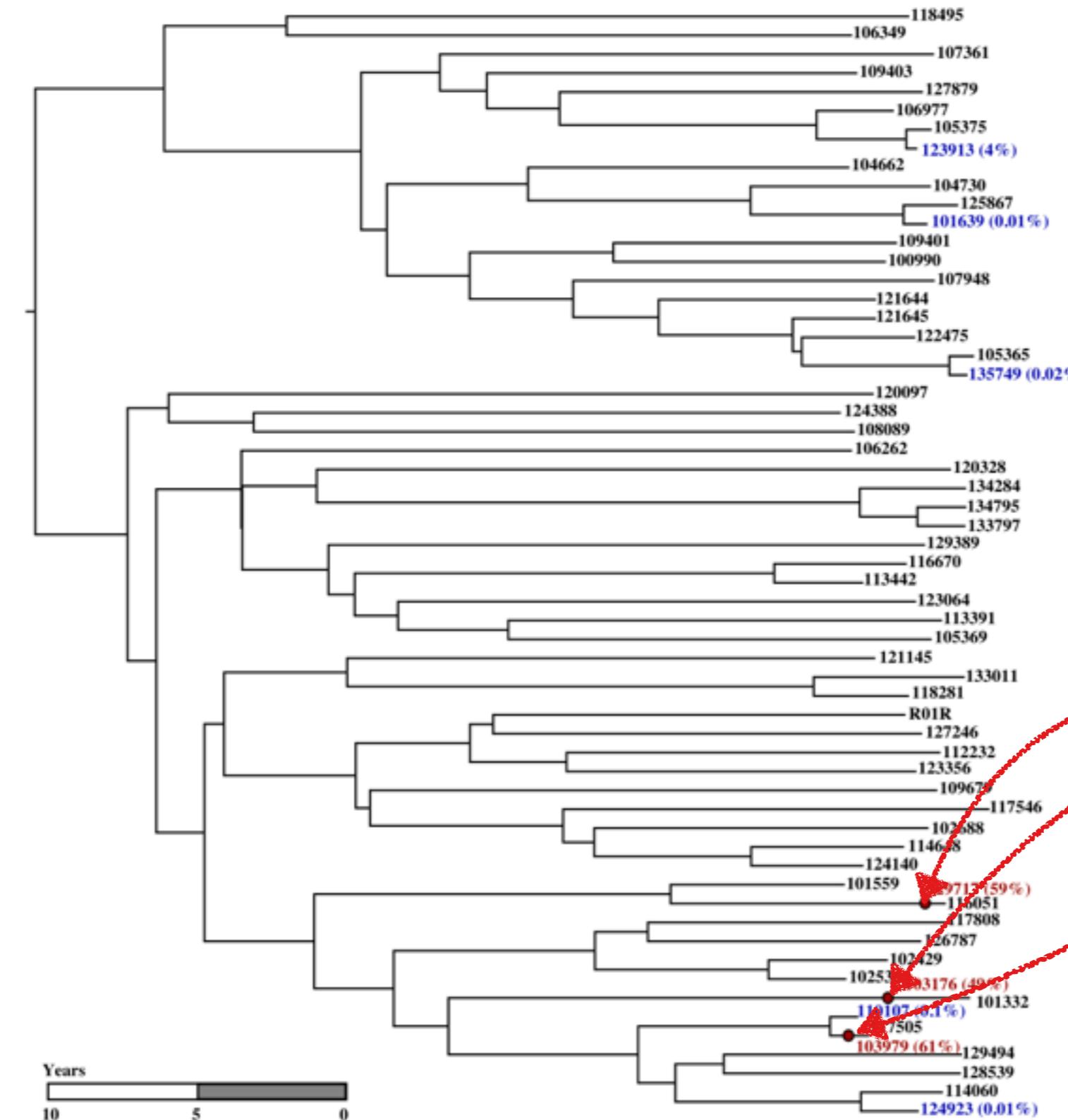
$\lambda_2, \mu_2, \psi_2, r_2$

$Q_2$

F



Sampled  
ancestors!



## Birth-death

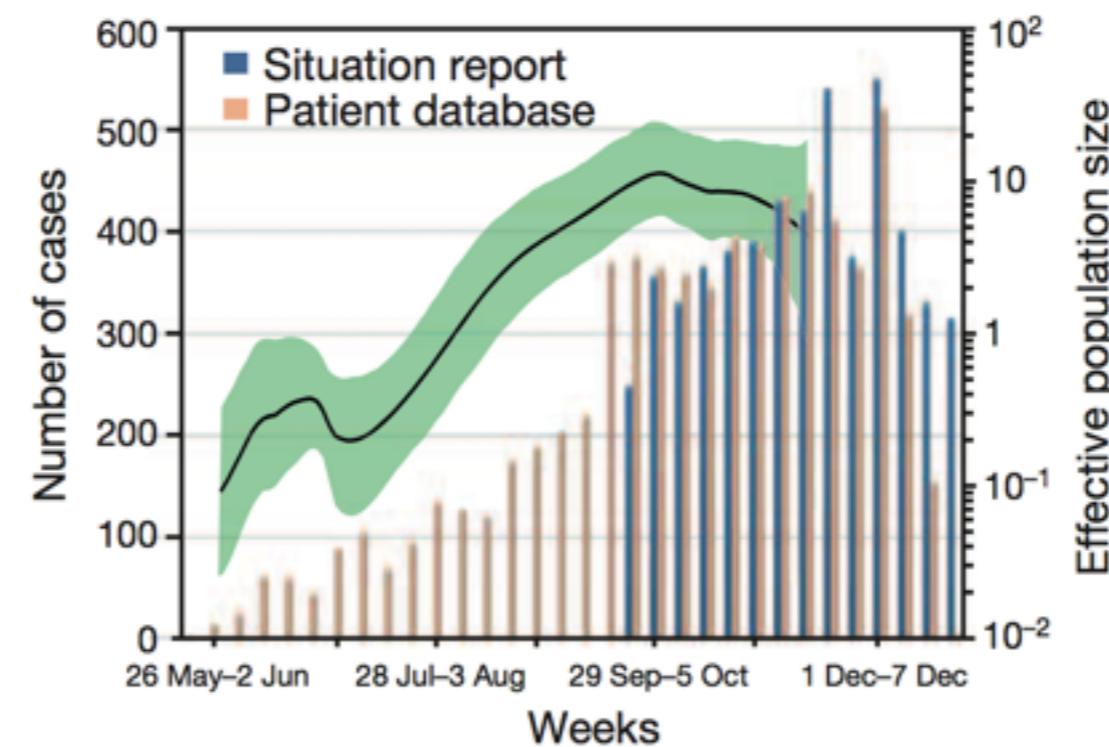
- Better treatment of stochasticity 😊
- Good for large sampling proportion and small population sizes 😊
- More computationally expensive 🤯
- Cannot easily calculate the likelihood for complex models 😢

## Coalescent

- Not as good at modelling stochasticity 😞
- Assumes a small sampling proportion 🤯
- Computationally relatively fast 😊
- Easier to derive coalescent times under complex models 😊

**But both models should give similar results most of the time...**

## Section 2: Parameter inference



Louis du Plessis

# Reminder: 2-step method for phylodynamic inference

---



## Step 1: Infer sampled transmission tree

$$P(\equiv | \mathcal{E}, \text{grid}, \text{clock})$$

Find model parameters that best describe the evolution of sequences along a tree

## Step 2: Infer dynamics that led to the tree

$$P(\mathcal{E} | \text{tree})$$

Find model parameters that best describe how the tree grows during an epidemic

**Tree and its dynamics are not always independent so a 1-step method would be better...**

$$P(\equiv | \mathcal{E}, \text{grid}, \text{clock}) P(\mathcal{E} | \text{tree})$$

## Maximum likelihood

- Find  $\theta$  that maximises:  
 $P(\mathcal{E} | \theta)$
- Point estimates
- May get stuck in local optima

## MCMC

- Find distribution of  $\theta$  given  $\mathcal{E}$
- Posterior of parameters:  
 $P(\theta | \mathcal{E})$
- Explore entire likelihood surface

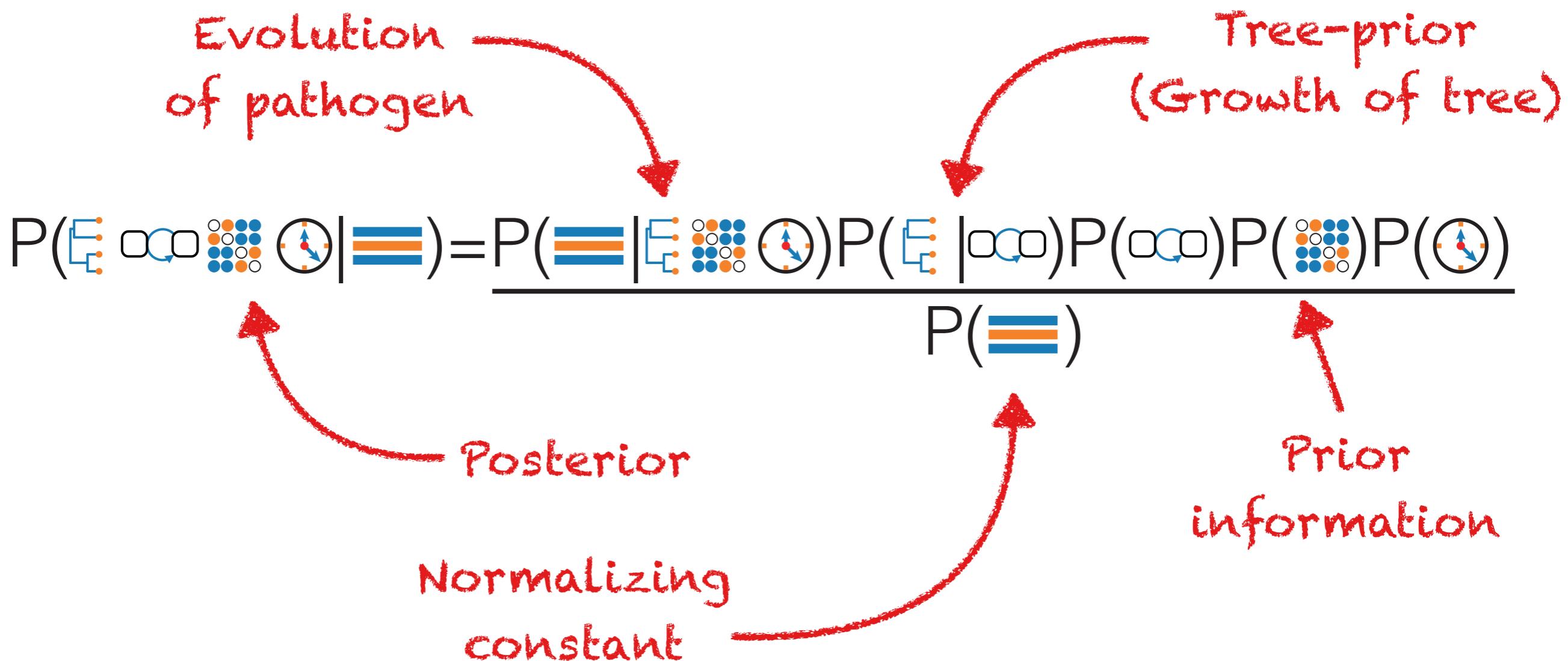
## Prior information

- Can include prior information in MCMC:

$$P(\theta | \mathcal{D}, \pi) = P(\theta)P(\mathcal{D} | \theta)\pi(\theta)$$

- Priors can bias results if too strong

# Posterior for MCMC inference



# 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

Epidemiological  
parameters directly  
from sequencing  
data!

$$P(E | \text{geno} \circ \text{demog} \circ \text{sub} \circ \text{clock}) = \frac{P(\text{geno} | E) P(\text{demog} | E) P(\text{sub} | E) P(\text{clock} | E)}{P(\text{geno})}$$

  
genetic  
sequences

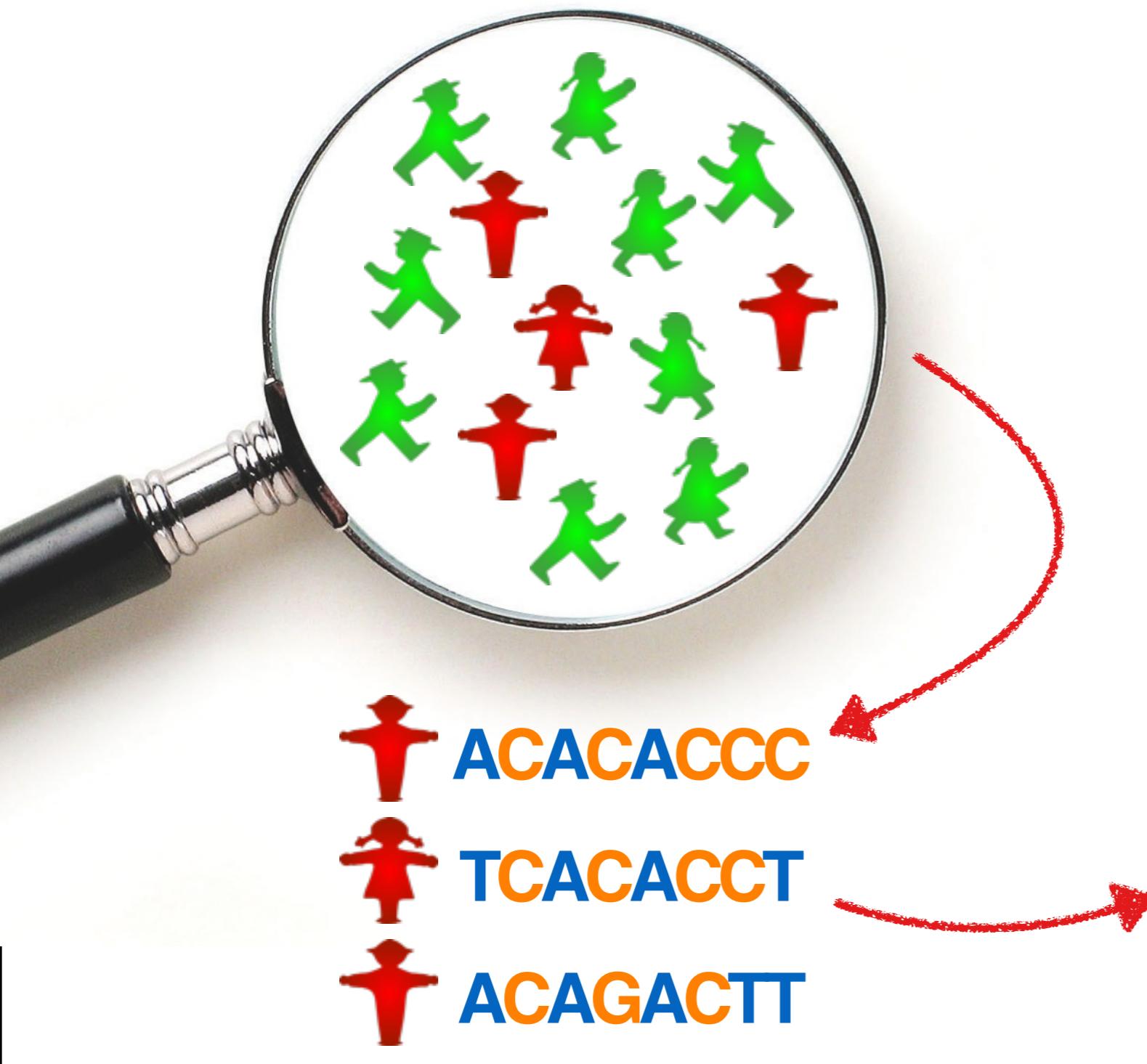
  
genealogy

  
demographic  
model

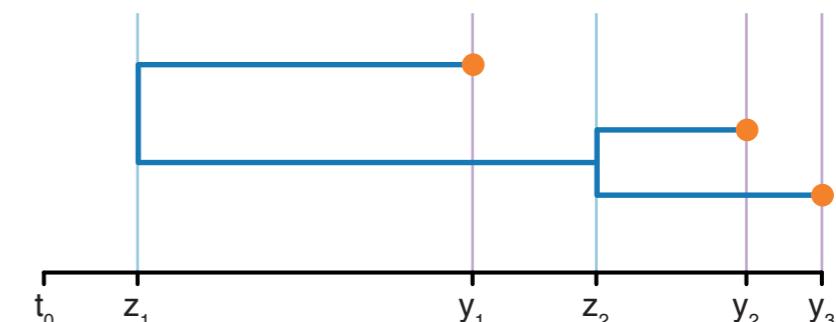
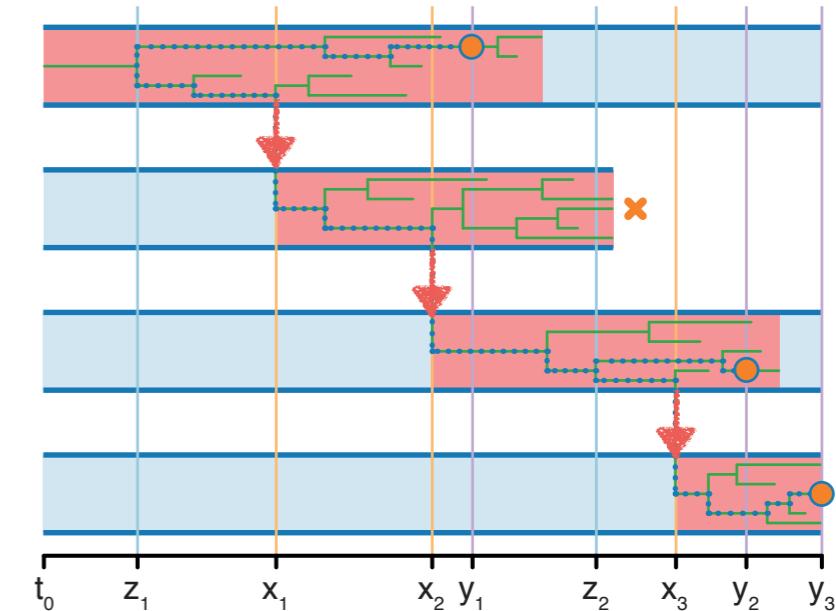
  
substitution  
model

  
molecular clock  
model

# Sampled transmission trees revisited



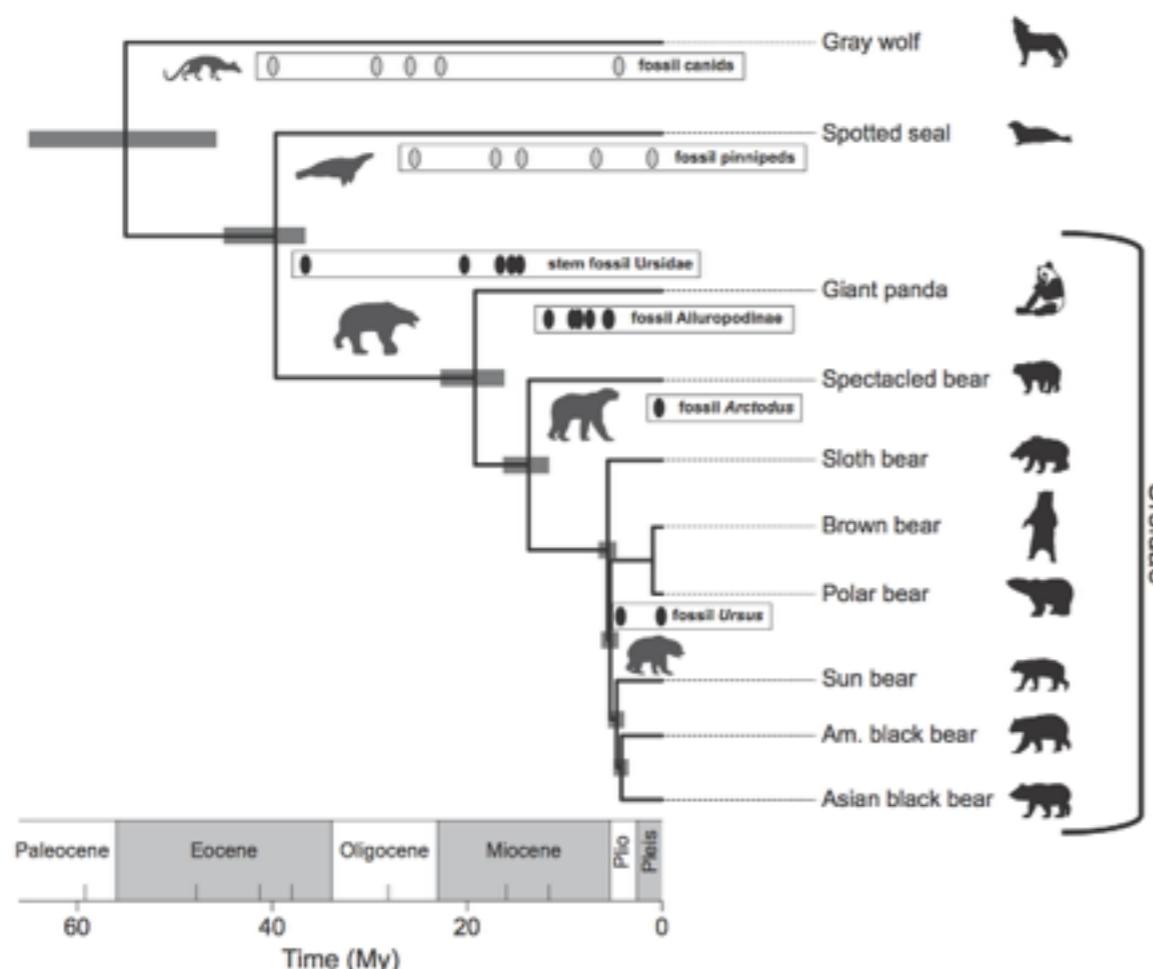
- No information on who infected whom
- Transmission events usually more recent than the MRCAs of the lineages



# Phylogenetics is not just about viruses!

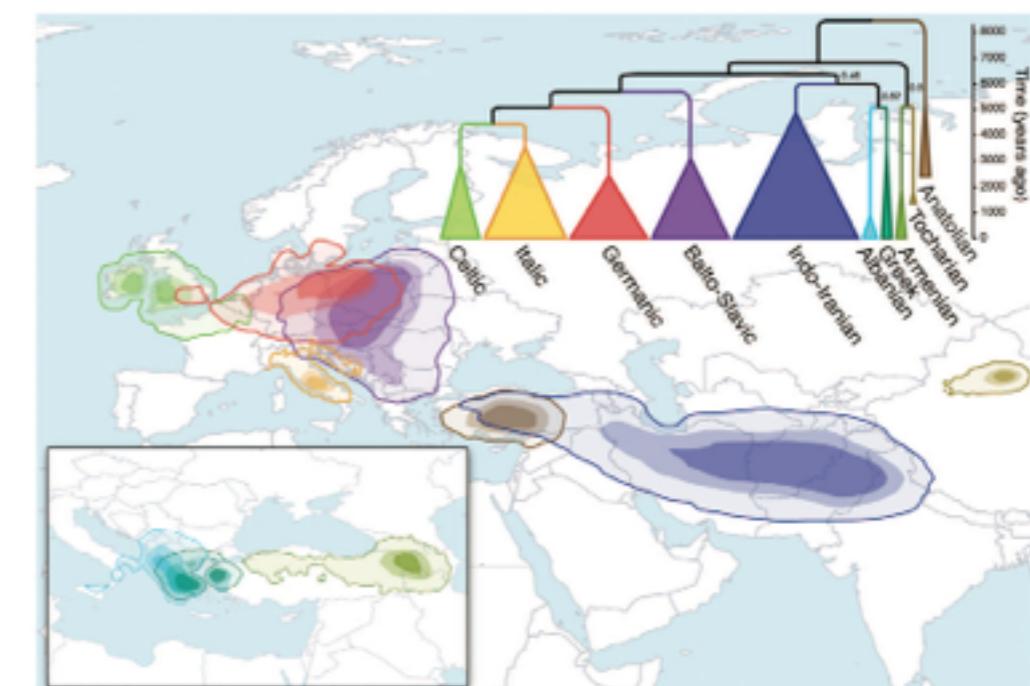
Phylogenetic methods are applicable to any infectious disease with enough phylogenetic diversity...

...or macro-evolution...



Heath et al. PNAS 2014

...or language evolution...

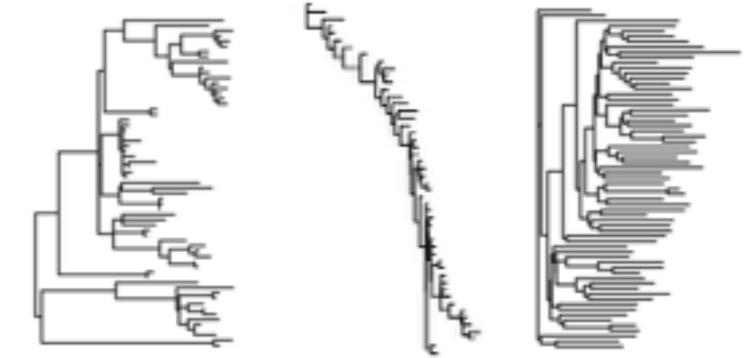


Bouckaert et al. Science 2012

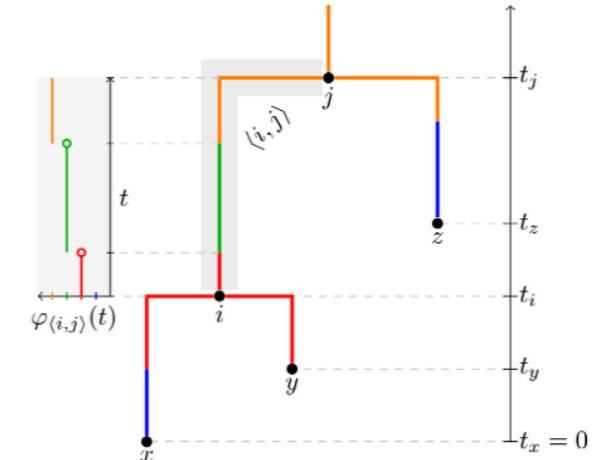
...and more!

# Summary

- Epidemiological dynamics leaves a signal on the branching pattern in the phylogeny
- Phylodynamic models can extract this signal by calculating the probability of the observed phylogeny conditioned on a demographic model
- Models can be specified in a birth-death or coalescent framework and may be more or less complex based on the dynamics that are to be modelled
- Using MCMC we can infer the posterior distributions of epidemiological parameters while accounting for phylogenetic uncertainty



$$P(\text{phylogeny} | \text{demographic model})$$



$$P(E|D) = \frac{P(D|E)P(E)P(D)P(\Theta)}{P(D)}$$