

Phylogenetic insights into infectious disease epidemiology

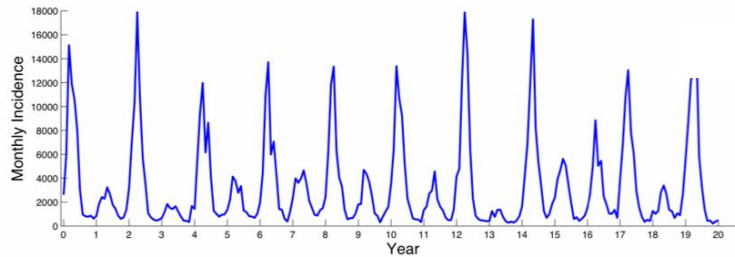
Molecular Epidemiology of Infectious Diseases
Lecture 1

January 10th, 2022

**Genomic data has
given us new power
to track the spread of
infectious pathogens**

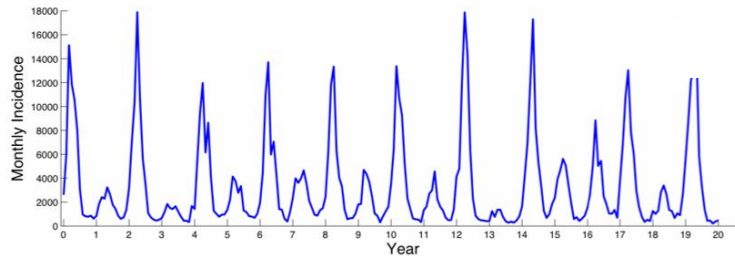
Revealing the source of infections

Classic sources of epidemiological data are typically not informative about the sources of new infections



Revealing the source of infections

Classic sources of epidemiological data are typically not informative about the sources of new infections



The genetic relatedness of pathogens sampled from different hosts or environments provides us with information about possible transmission routes including **the source of new infections.**

Hourglass format of course



Starting from very different backgrounds

Core phylogenetic methods applicable across systems

More targeted applications and team projects

Coursework and grades

“Everyone should get a A”

There is no graded work other than a team project focusing on a pathogen and dataset of your choice during the second half of the semester.

But please do:

- Look at the suggested readings.
- Participate in in-class discussions and tutorials
- Come to class ready to ask questions and discuss problems

The importance of phylogenies

While there are many methods for analyzing pathogen genomic data, this lecture and most of this class will examine phylogenetic methods.

Phylogenies describe the ancestral relationships among individuals or taxa in terms of shared descent.

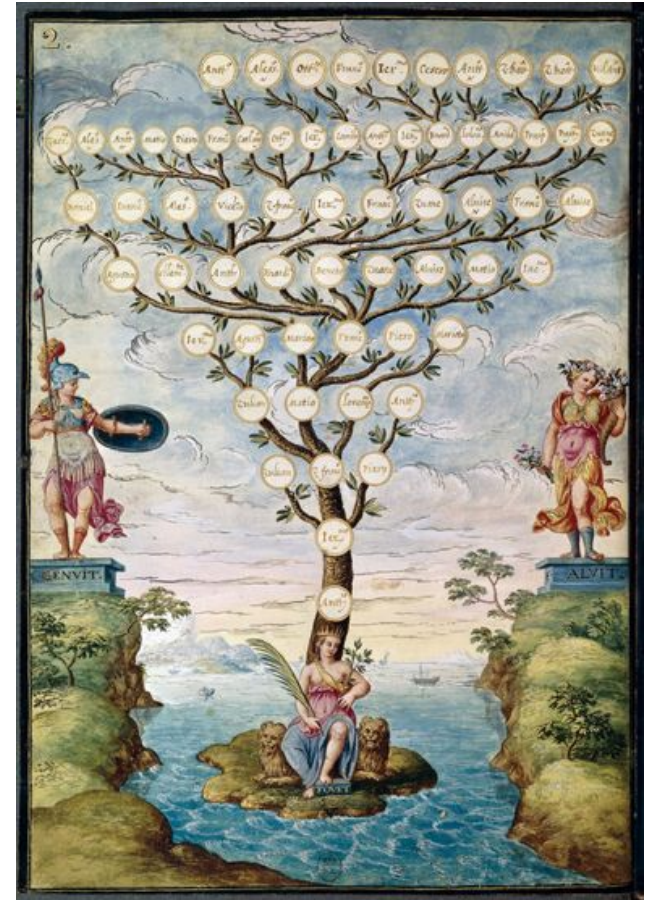


Image from *The Book of Trees* (Manuel Lima, 2014)

Why phylogenies?

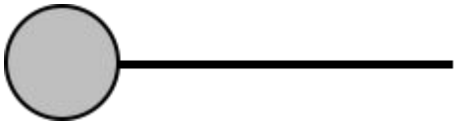
1. The branching structure of trees can be directly related back to the transmission dynamics of a pathogen
2. Thinking phylogenetically can help us understand how epidemic dynamics shape genetic variation in a pathogen population.



Image from *The Book of Trees* (Manuel Lima, 2014)

Let's start by
considering a small
epidemic spreading
through a host
population

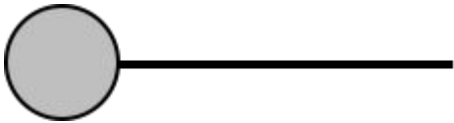
A simple epidemic example



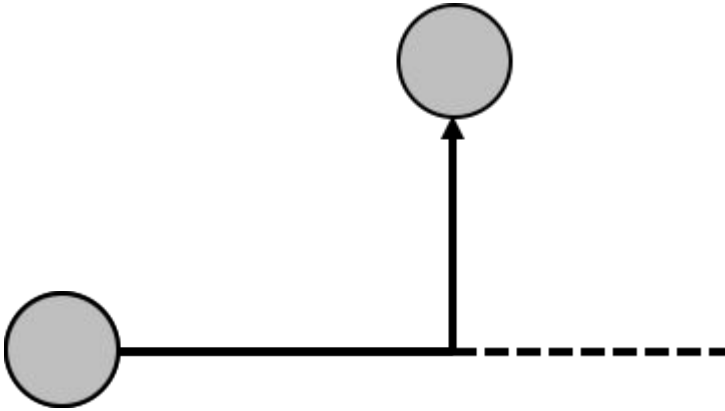
A simple epidemic example



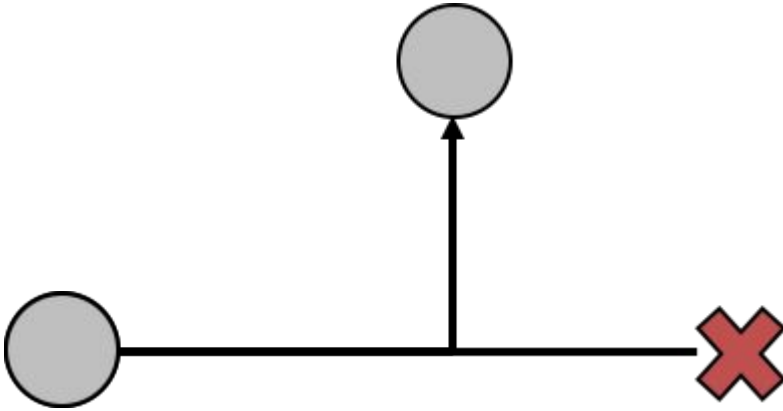
A simple epidemic example



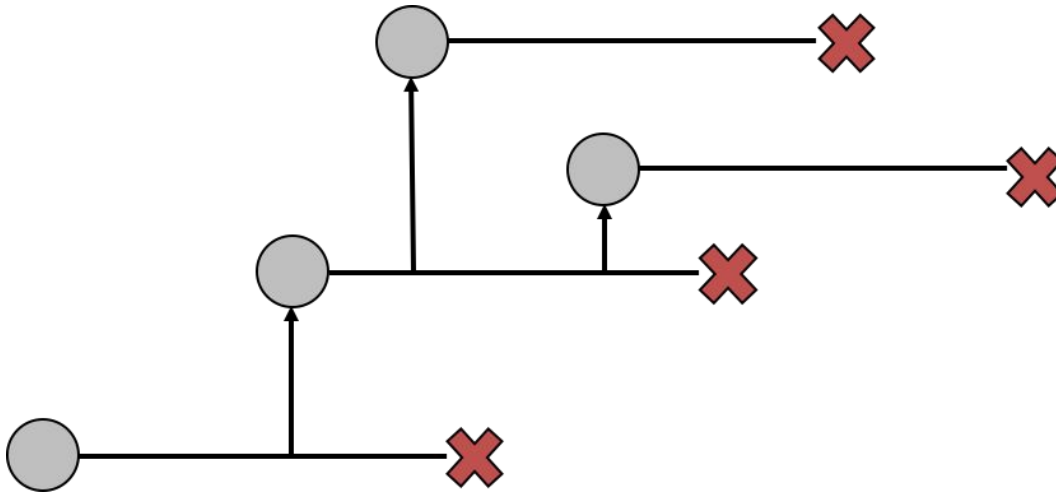
A simple epidemic example



A simple epidemic example

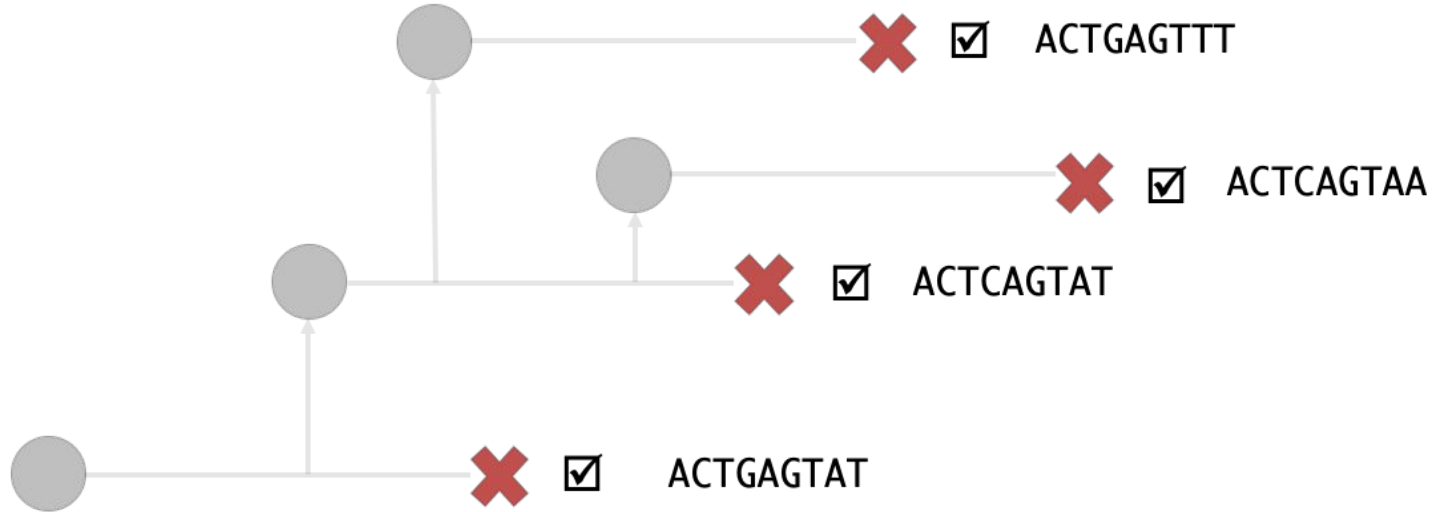


A simple epidemic example

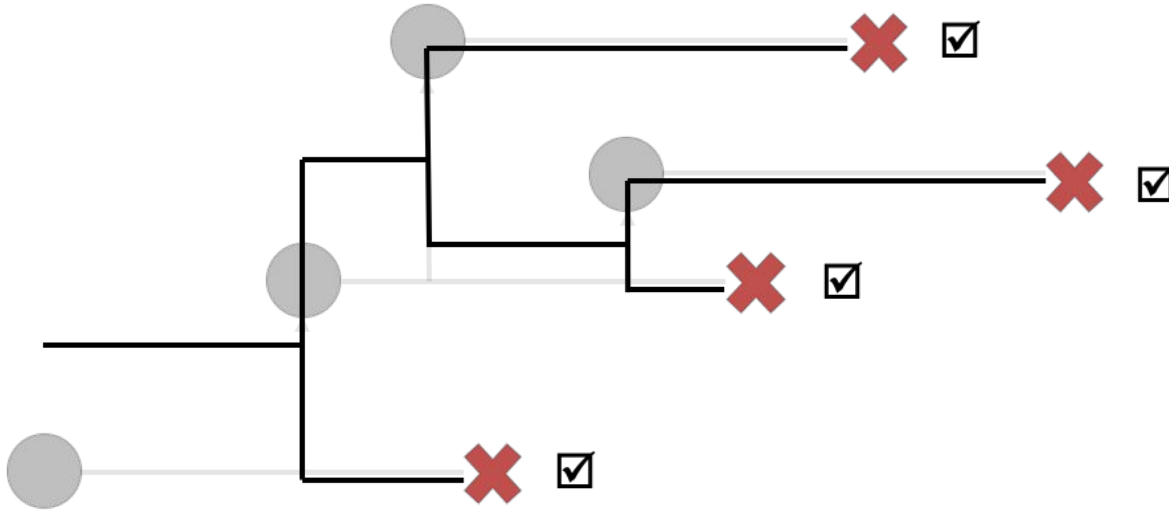


Transmission tree

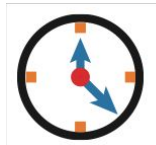
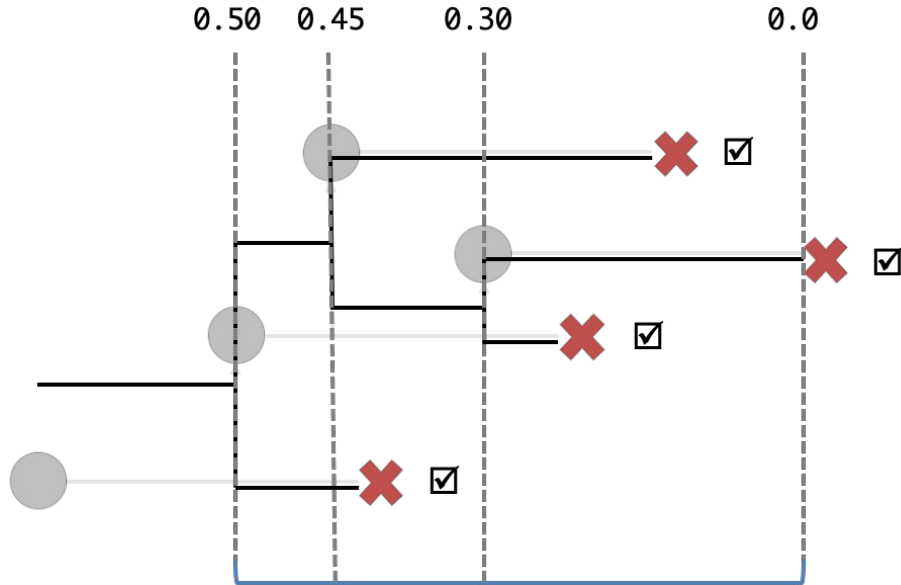
A simple epidemic example



A simple epidemic example



A simple epidemic example



$$\text{Real time} = \frac{\text{genetic distance}}{\text{clock rate}}$$

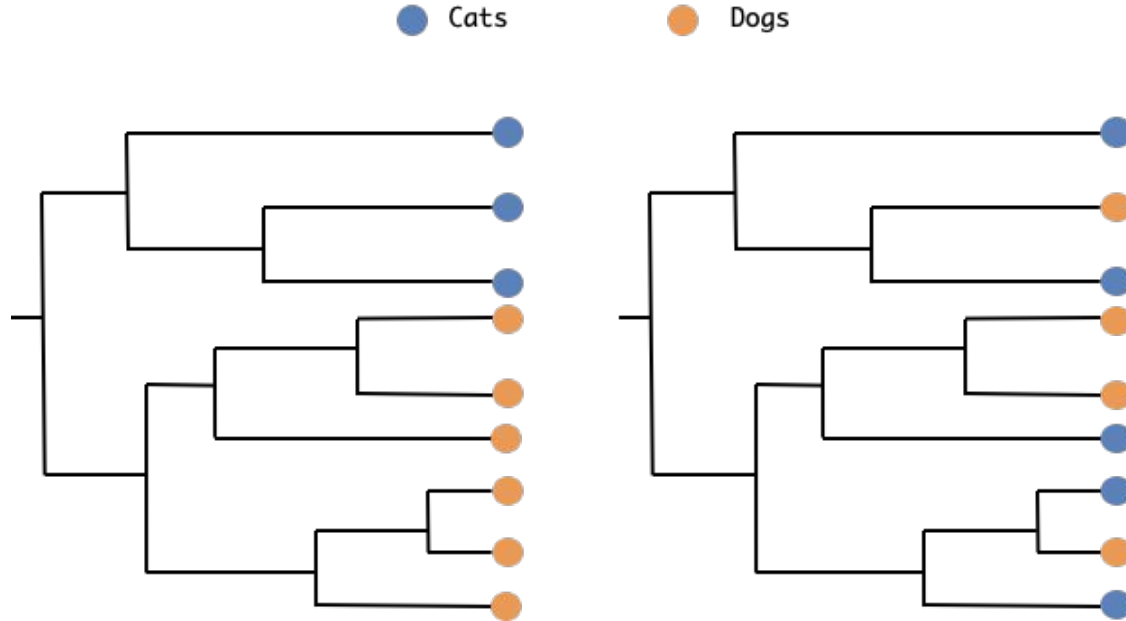
Phylogenies can tell us about:

- Linkage and the sources of transmission
- The origins of epidemics and new strains
- Past epidemic dynamics
- Pathogen fitness and adaptation

Phylogenies can tell us about:

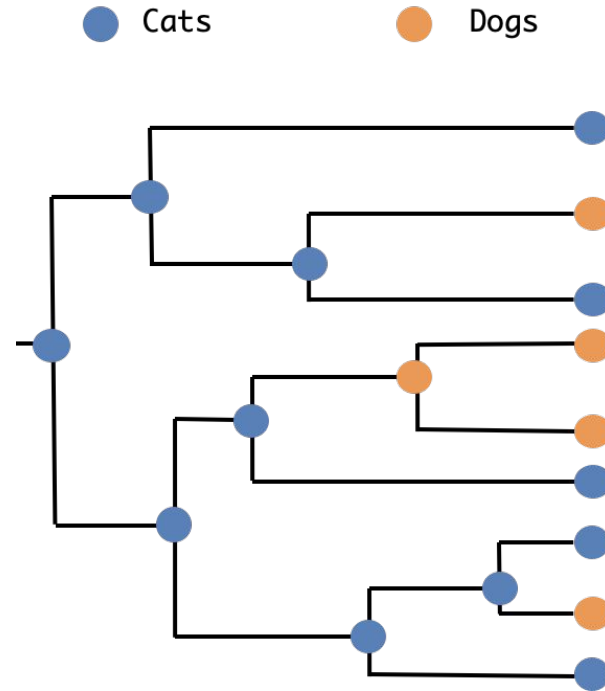
- Linkage and the sources of transmission
- The origins of epidemics and new strains
- Past epidemic dynamics
- Pathogen fitness and adaptation

Phylogenetic linkage



Ancestral state reconstruction

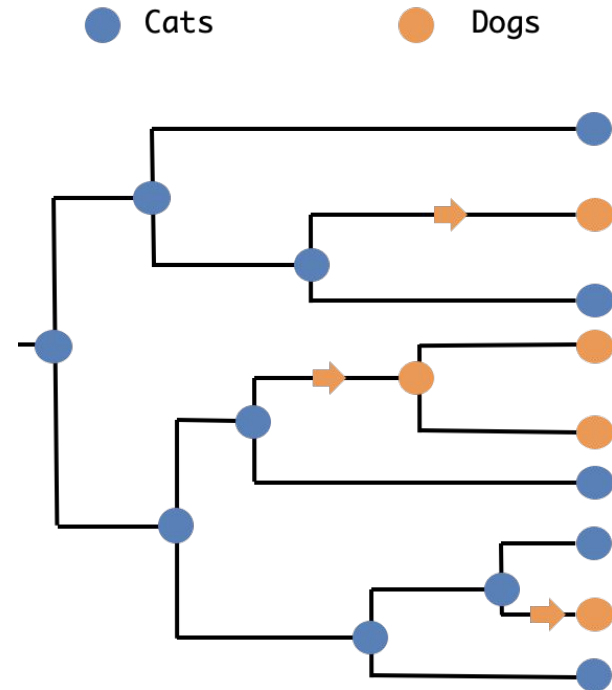
Ancestral state reconstruction allows us to infer the location/host of past transmission events.



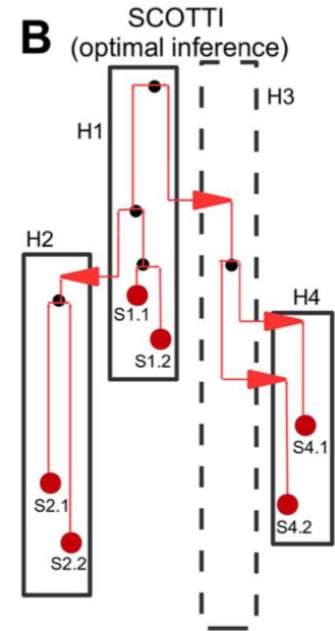
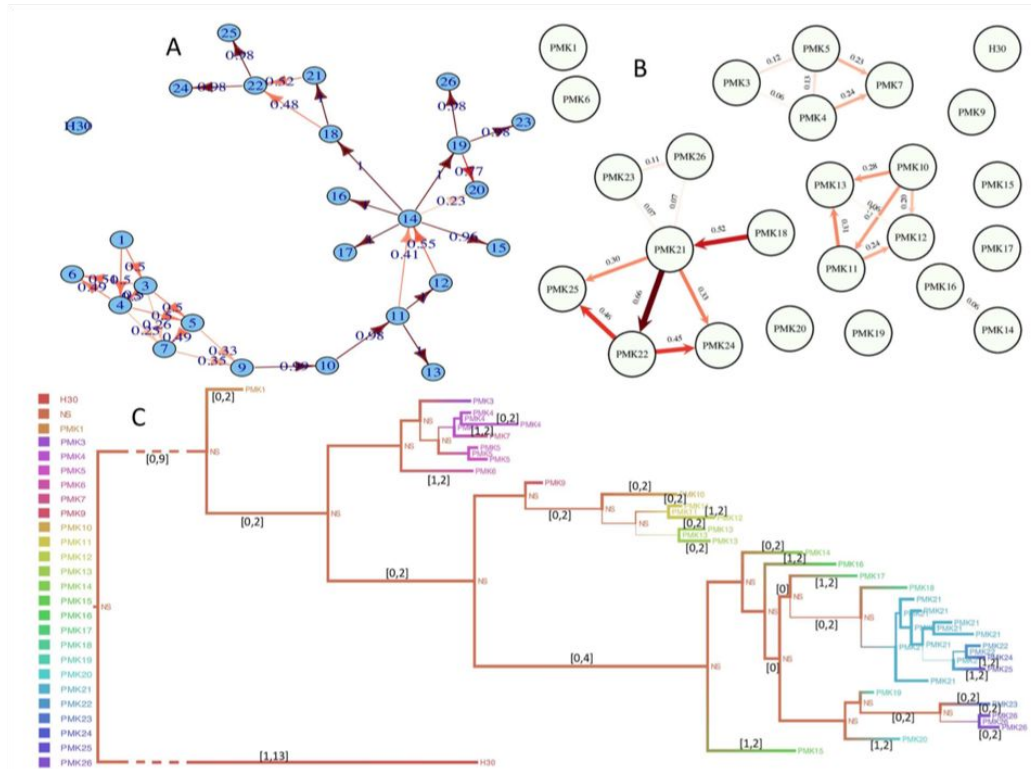
Ancestral state reconstruction

Ancestral state reconstruction allows us to infer the location/host of past transmission events.

Ancestral states can therefore allow us to infer the direction of infection.



Klebsiella transmission trees

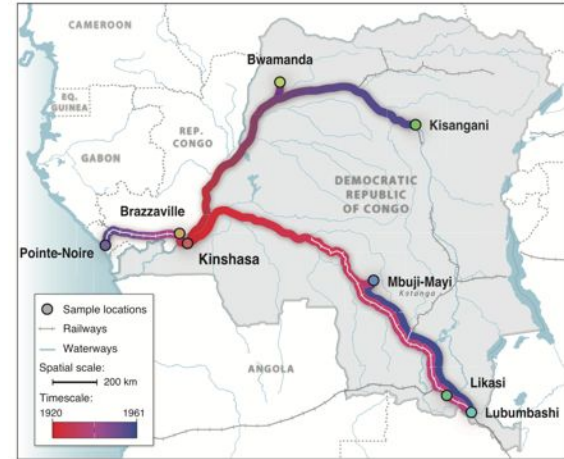
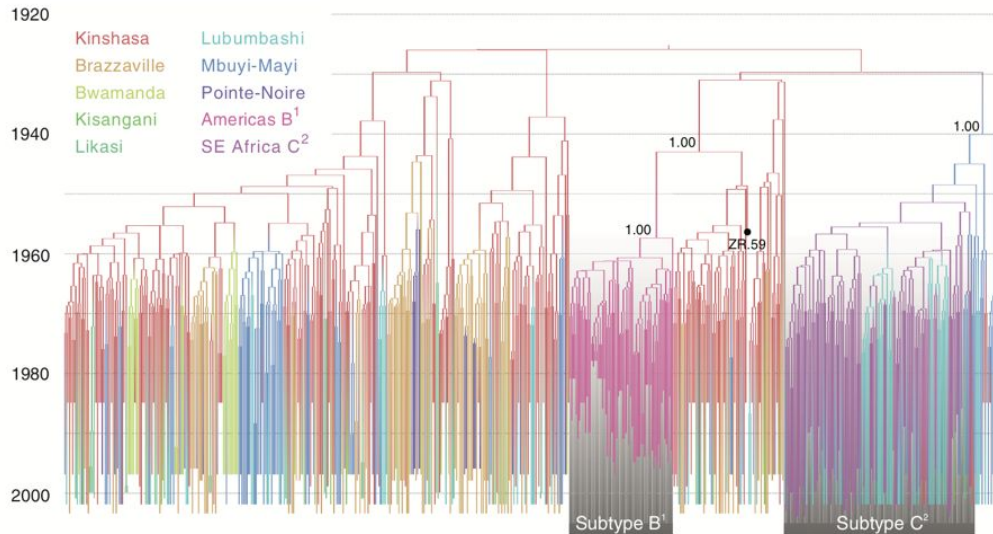
De Maio *et al.* (PCB, 2016)

Phylogenies can tell us about:

- Linkage and the sources of transmission
- The origins of epidemics and new strains
- Past epidemic dynamics
- Pathogen fitness and adaptation

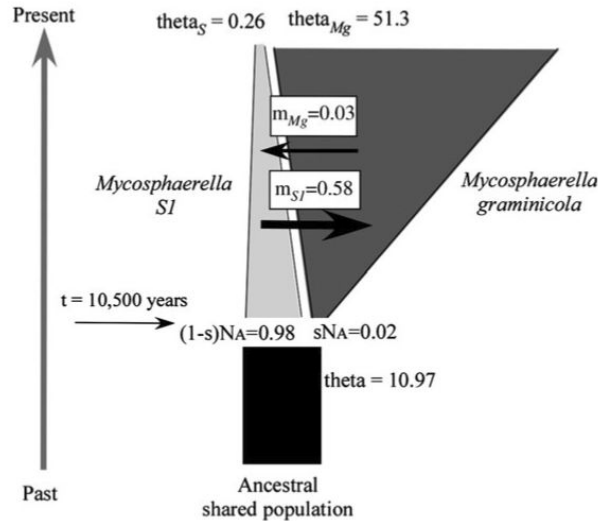
Origins of the HIV-1 epidemic

Faria *et al.* (Science, 2014) traced the origins of the HIV-1 epidemic back to the 1920's and 30's in Kinshasa, DRC.



Origins of *Mycosphaerella graminicola*

Stukenbrock *et al.* (MBE, 2006) traced the fungal pathogen causing septoria leaf blotch on wheat back to 8,000 to 9,000 BC in the Fertile Crescent.



M. graminicola on wheat (Wikipedia)

Neolithic origins of other agro-pathogens

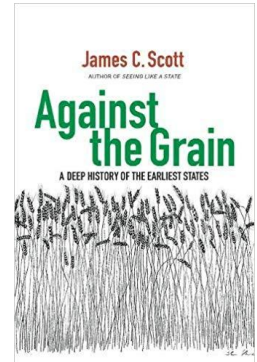
Supports idea that many agriculturally important pathogens arose during the Neolithic transition to farming.

Table 1 Examples of evolutionary mechanisms by which plant pathogens have emerged in agro-ecosystems over different time scales

Evolutionary mechanism	Plant pathosystem	Time scale	Reference
Domestication/host-tracking			
	<i>Mycosphaerella graminicola</i> on wheat	10–12,000 years BP	95
	<i>Magnaporthe oryzae</i> on rice	7000 years BP	24
	<i>Phytophthora infestans</i> on potato	7000 years BP	34
	<i>Ustilago maydis</i> on maize	8000 years BP	72
Host jump/host shift			
	<i>Magnaporthe oryzae</i> from <i>Setaria</i> millet to rice	Abrupt evolutionary change, approx. 7000 years BP	24
	<i>Rhynchosporium secalis</i> from wild grasses to barley and rye	Abrupt evolutionary change, approx. 2,000 years BP	111
	<i>Phytophthora infestans</i> from wild <i>Solanum</i> species to potato	Abrupt evolutionary change, <500 years BP	35, 39

Stukenbrock and McDonald (Annu. Rev. Phyto., 2008)

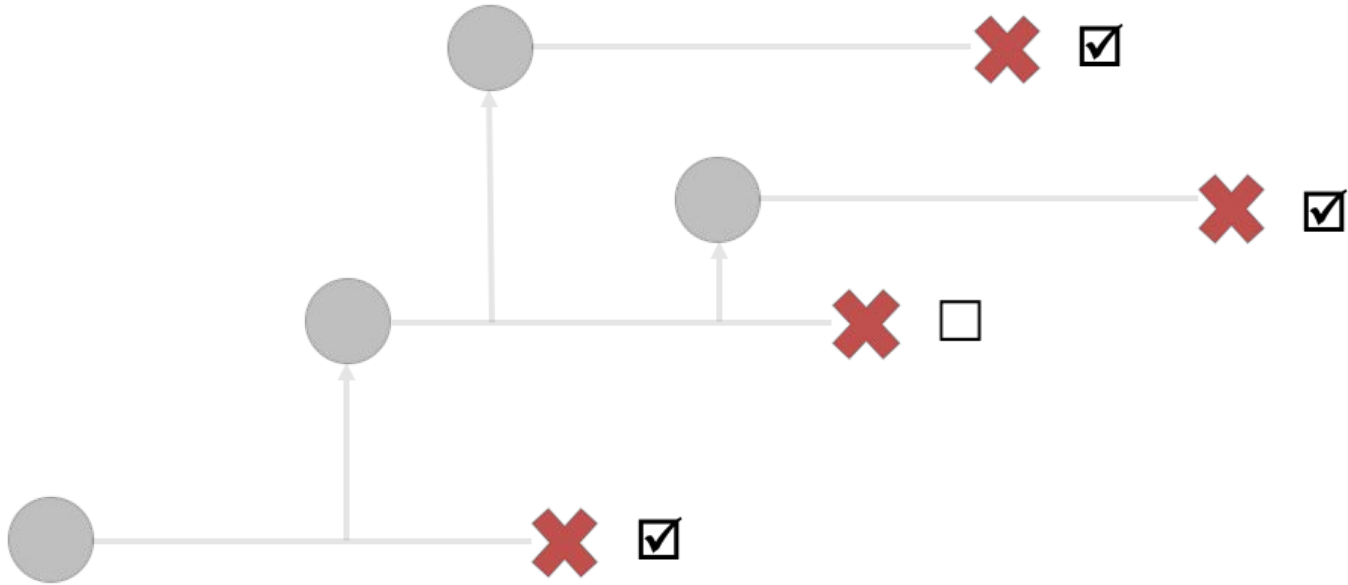
“Neolithic pathogen relocation camps”



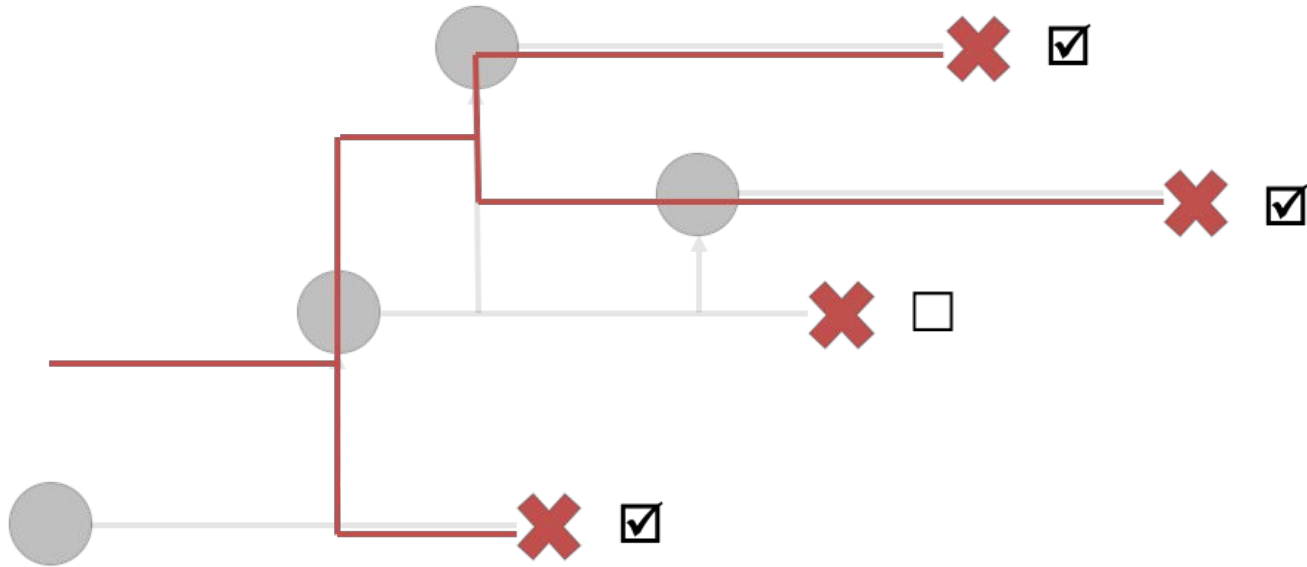
Phylogenies can tell us about:

- Linkage and the sources of transmission
- The origins of epidemics and new strains
- Past epidemic dynamics
- Pathogen fitness and adaptation

A simple epidemic example with incomplete sampling



A simple epidemic example with incomplete sampling



We only observe transmission events as branching events if we sample both the parent and child lineage descending from the transmission event

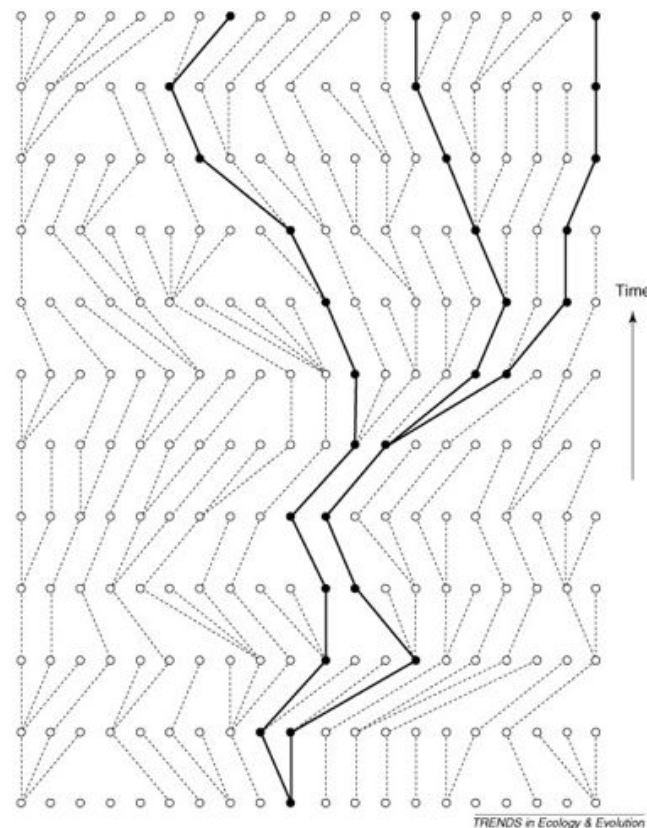
This brings us to
phylodynamic
modeling

Phylogenetic modeling in a nutshell

Phylogenies will only contain sampled lineages.

The sampled lineages are embedded within the full ancestral history of the population.

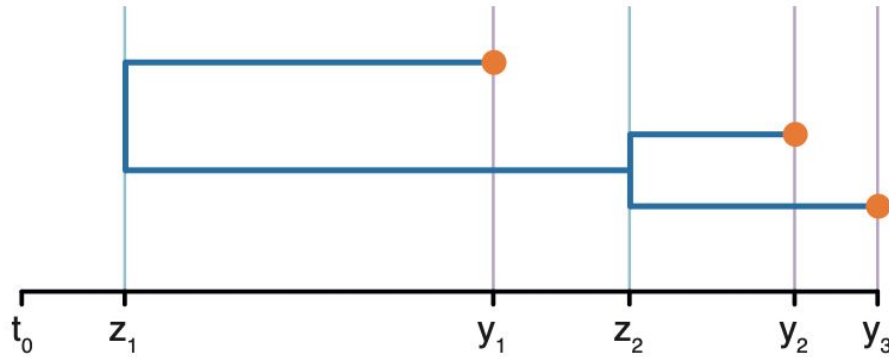
We need a statistical model that allows us to infer the most likely population history from the sampled phylogeny.



Kuhner *et al.* (2008)

Two types of phylodynamic models

Birth-death →



← **Coalescent**

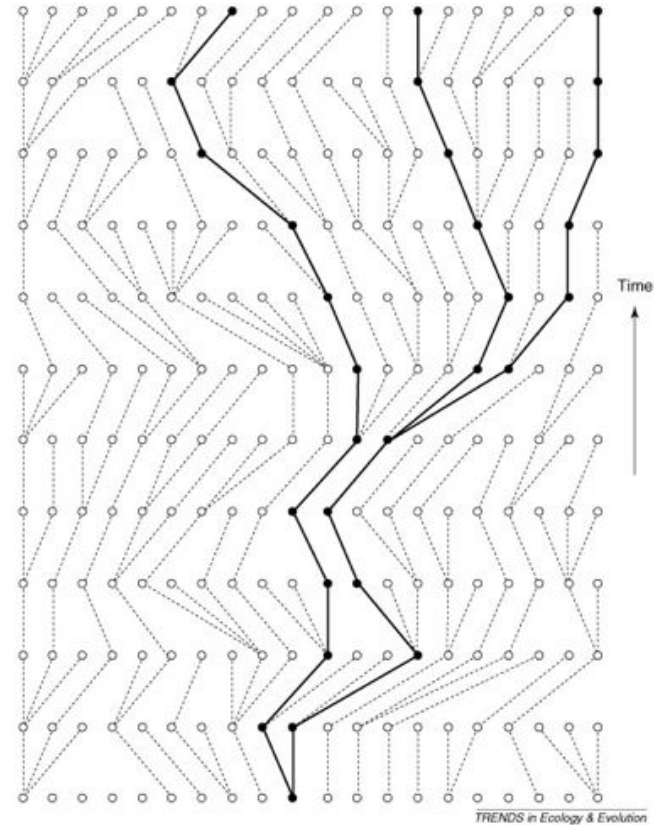
Coalescent theory

The coalescent traces the ancestry of sampled individuals back in time.

Allows us to relate events observed in the tree to the larger history of a population

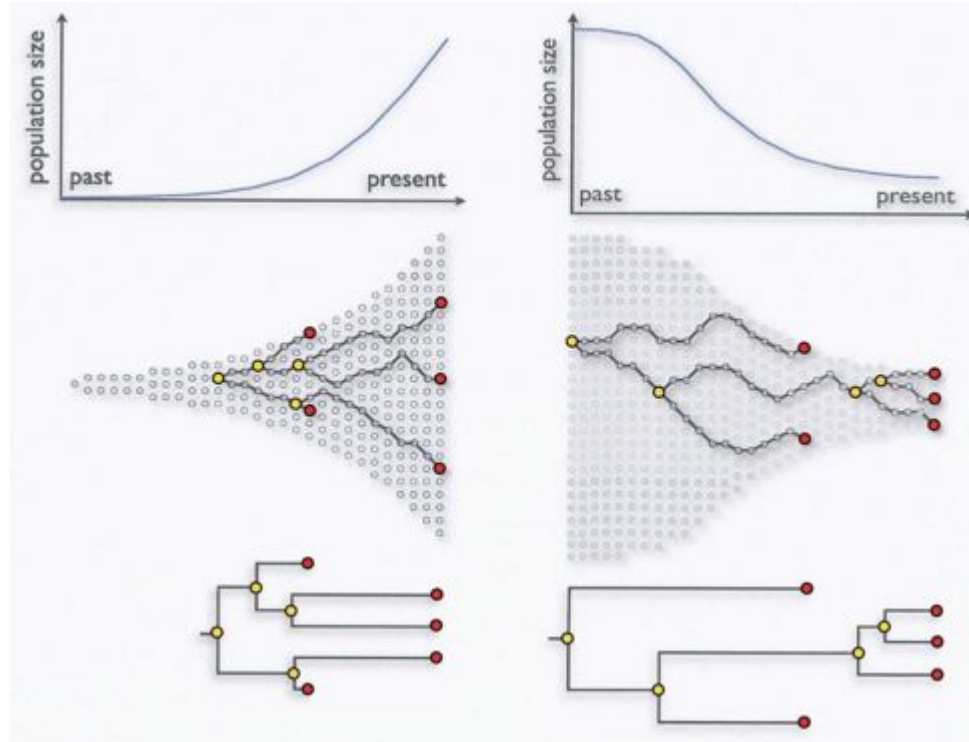
Probability of two lineages coalescing per generation is:

$$p_{coal} = \frac{1}{N}$$



Kuhner *et al.* (2008)

Reconstructing population dynamics

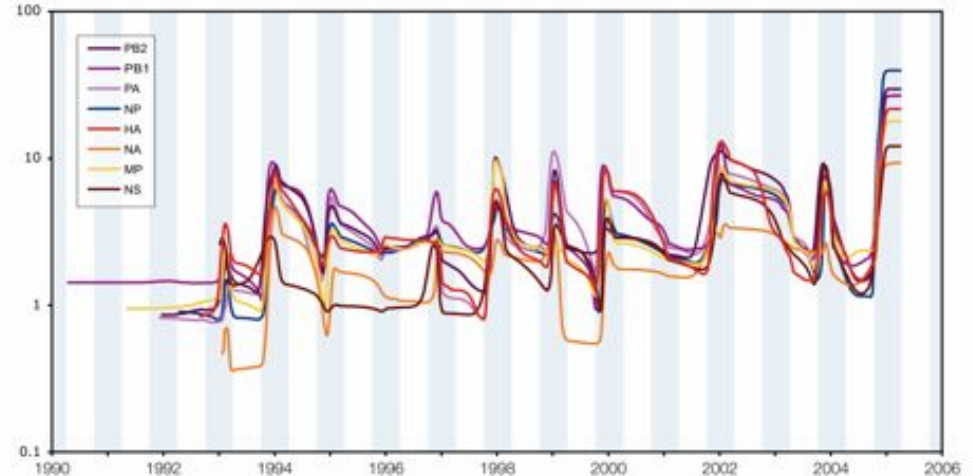


Reconstructing dynamics: influenza A

The genomic and epidemiological dynamics of human influenza A virus

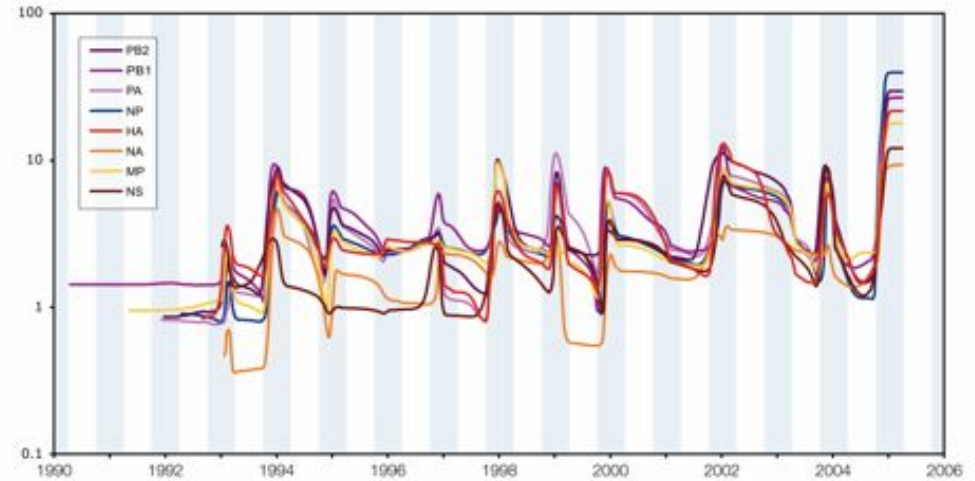
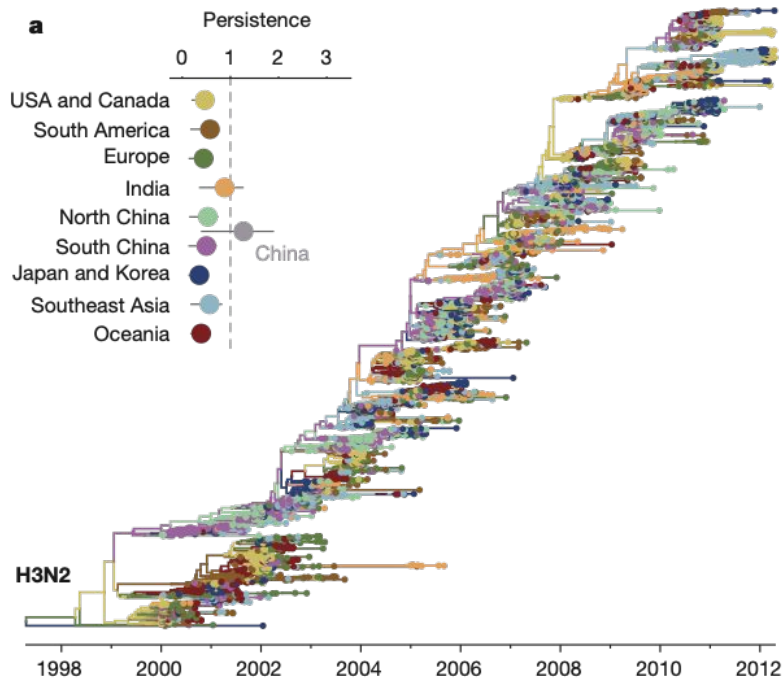
Andrew Rambaut¹, Oliver G. Pybus², Martha I. Nelson³, Cecile Viboud¹, Jeffery K. Taubenberger³ & Edward C. Holmes^{1,4}

The evolutionary interaction between influenza A virus and the human immune system, manifest as 'antigenic drift' of the viral haemagglutinin, is one of the best described patterns in molecular evolution. However, little is known about the genome-scale evolutionary dynamics of this pathogen. Similarly, how genomic processes relate to global influenza epidemiology, in which the A/H3N2 and A/H1N1 subtypes co-circulate, is poorly understood. Here through an analysis of 1,302 complete viral genomes sampled from temperate populations in both hemispheres, we show that the genomic evolution of influenza A virus is characterized by a complex interplay between frequent reassortment and periodic selective sweeps. The A/H3N2 and A/H1N1 subtypes exhibit different evolutionary dynamics, with diverse lineages circulating in A/H1N1, indicative of weaker antigenic drift. These results suggest a sink-source model of viral ecology in which new lineages are seeded from a persistent influenza reservoir, which we hypothesize to be located in the tropics, to sink populations in temperate regions.



Rambaut *et al.* (2008)

Reconstructing dynamics: influenza A

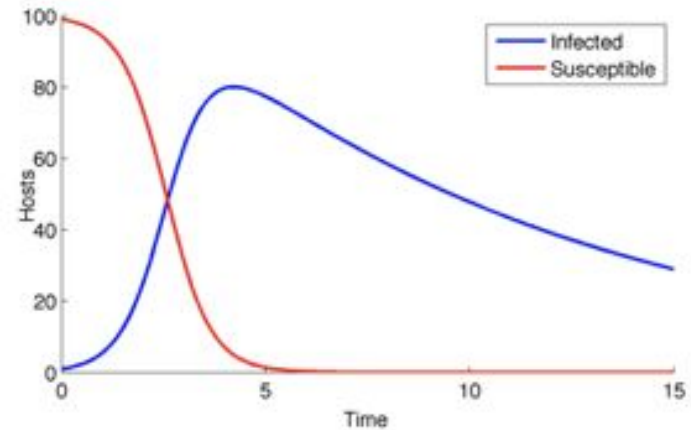
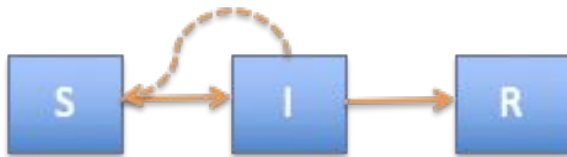


Rambaut *et al.* (2008)

Bedford *et al.* (Nature, 2015)

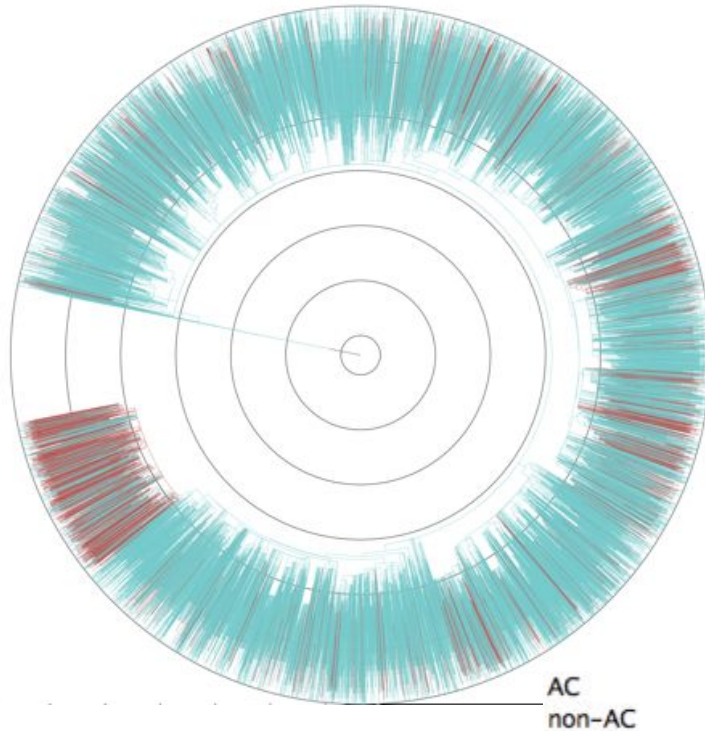
Coupling epidemiological models to trees

We can use phylodynamic modeling to couple phylogenetic methods with more traditional epidemiological models

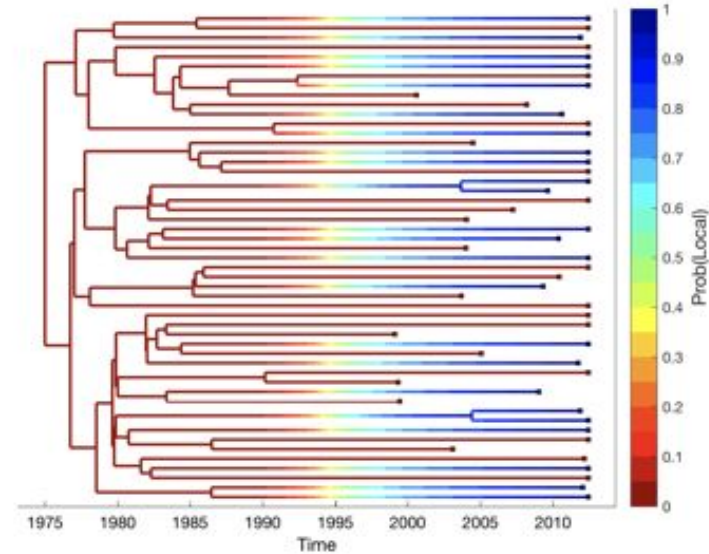
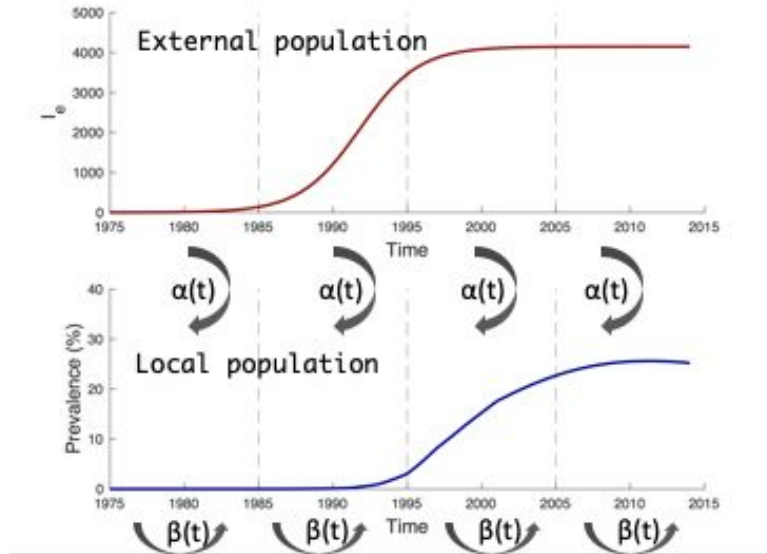


We can formulate epidemic models that we can then fit to phylogenies to estimate parameters of interest.

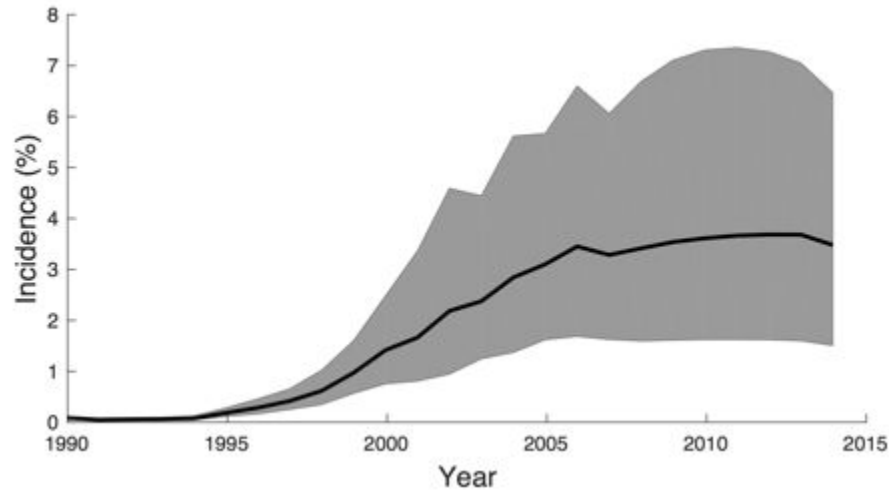
HIV in rural Kwa-Zulu Natal



A simple two-patch SIR model for HIV



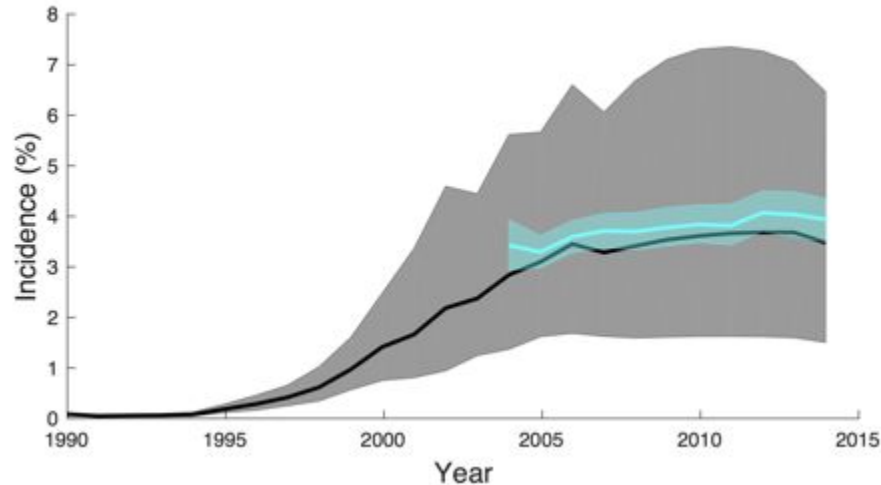
Phylodynamic estimates of HIV incidence



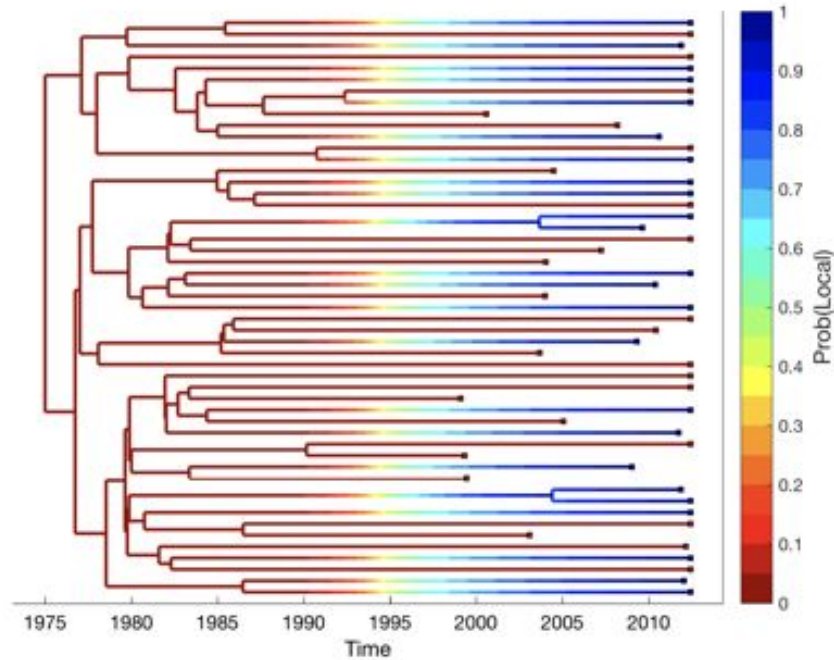
Rasmussen et al. (Virus Evolution, 2018)

Phylodynamic estimates of HIV incidence

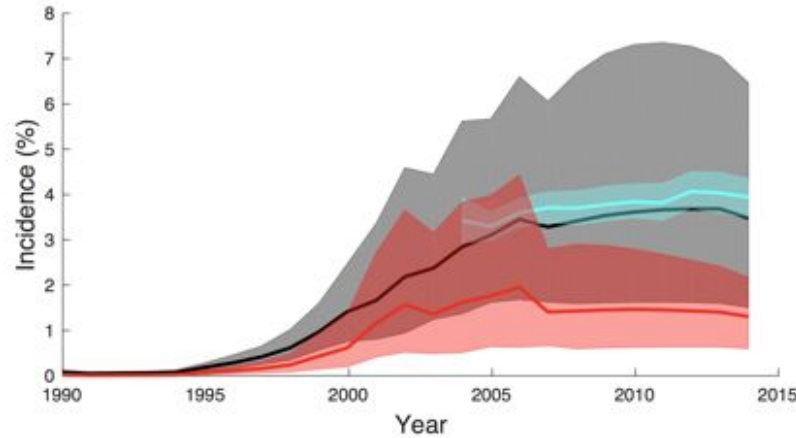
Inferred incidence of 3-4% per year almost perfectly coincides with population-based surveillance data.



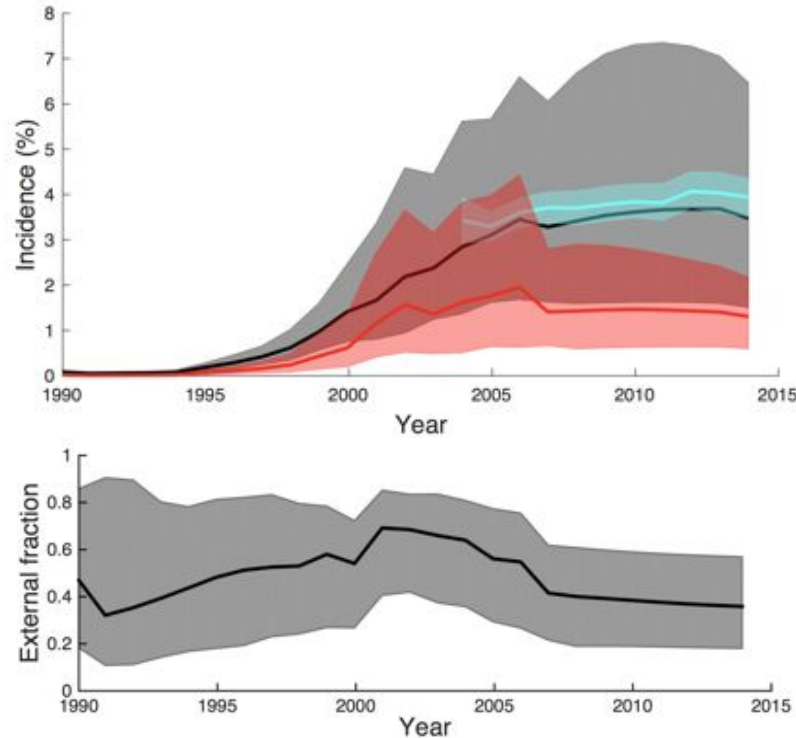
Tracking lineage movement



Incidence due to external introductions



Incidence due to external introductions



As of 2014, 35% of new infections were attributed to external introductions.

Phylogenies can tell us about:

- Linkage and the sources of transmission
- The origins of epidemics and new strains
- Past epidemic dynamics
- Pathogen fitness and adaptation

Phylodynamics with selection

Selection for higher fitness strains strongly shapes the phylogenetic history of many different pathogens.

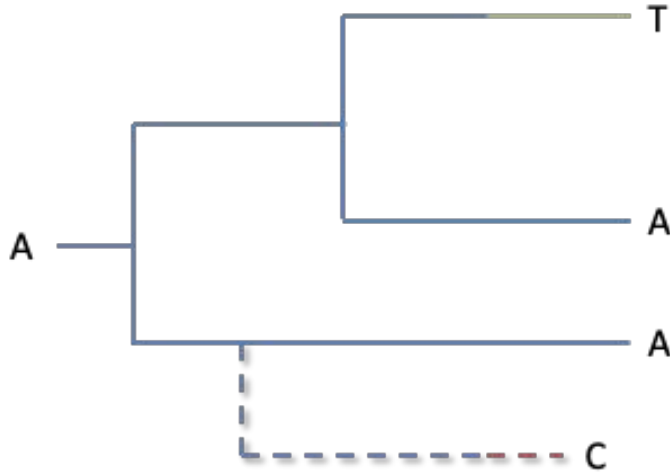


Grenfell *et al.* (Science, 2004)

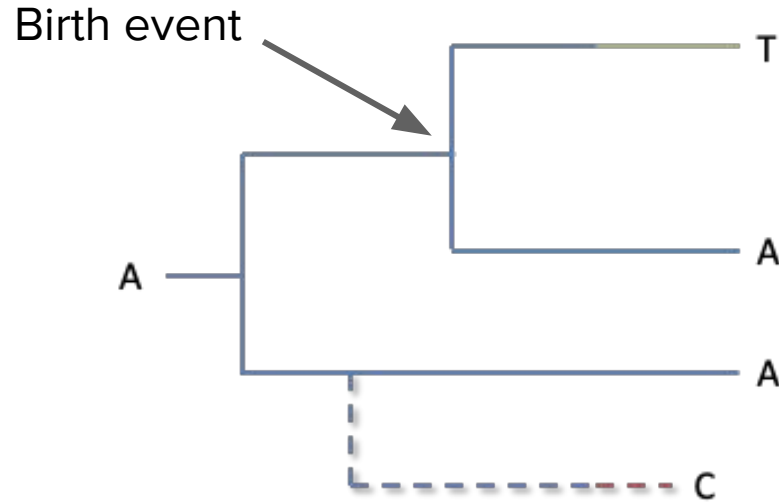
**We therefore need
phylodynamic models
that allow selection to
shape trees**

Multi-type birth-death models

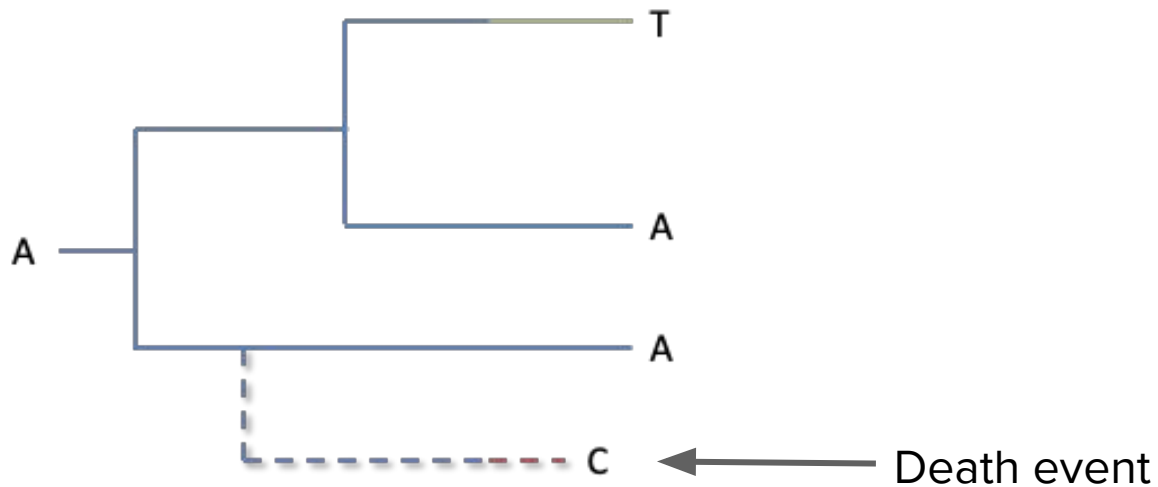
Provide one way of incorporating adaptive (non-neutral) evolution into phylogenetic models.



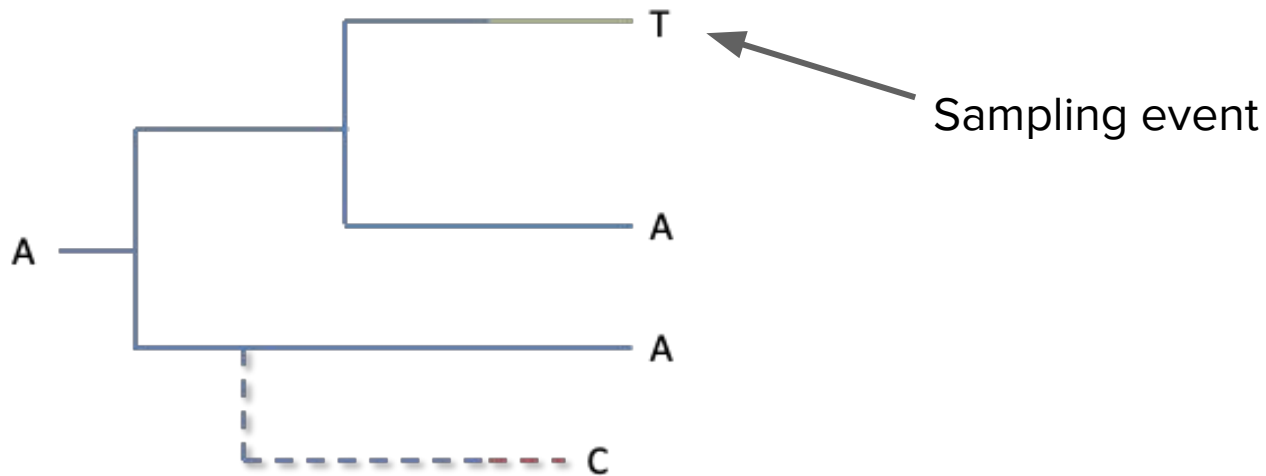
Multi-type birth-death models



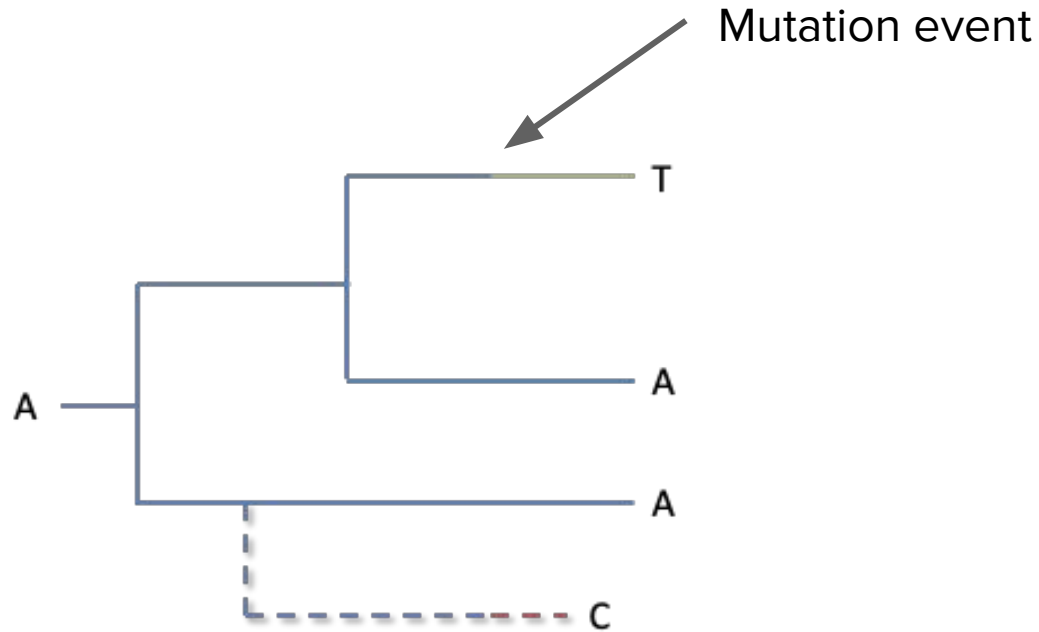
Multi-type birth-death models



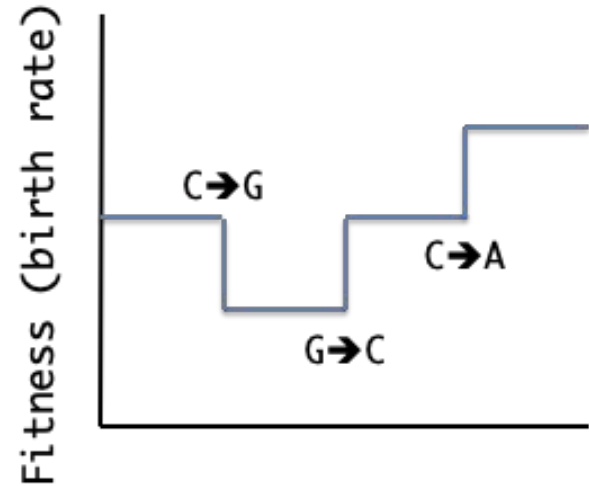
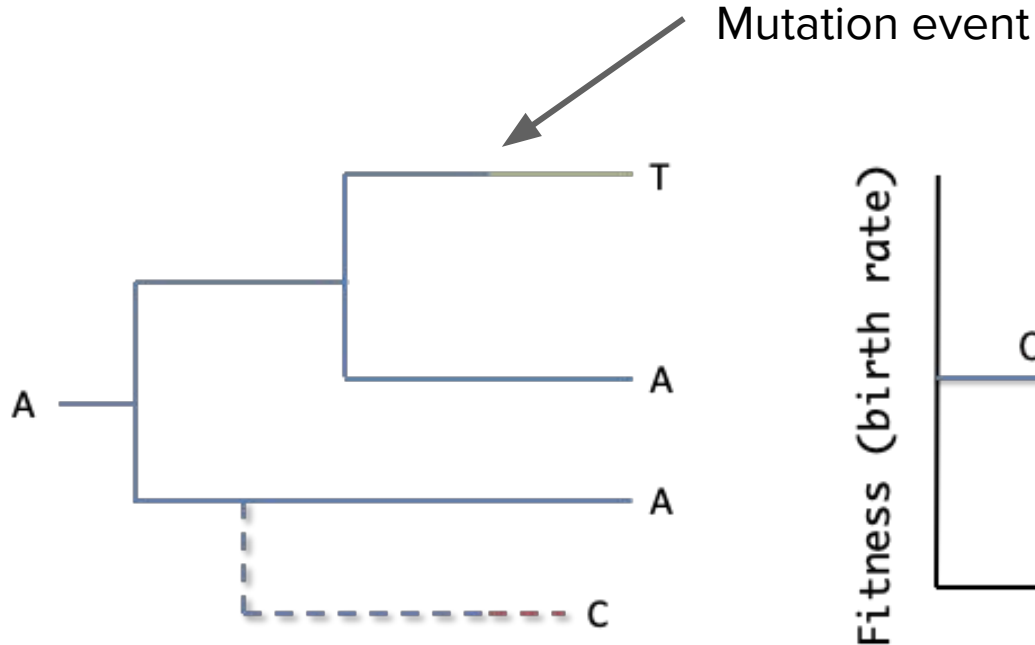
Multi-type birth-death models



Multi-type birth-death models

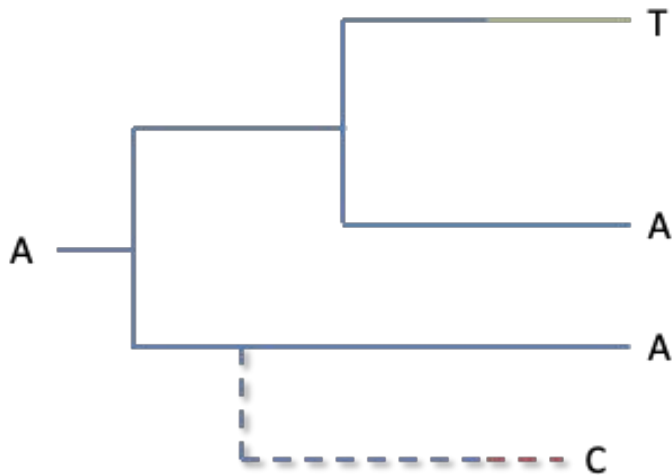


Multi-type birth-death models



Multi-type birth-death models

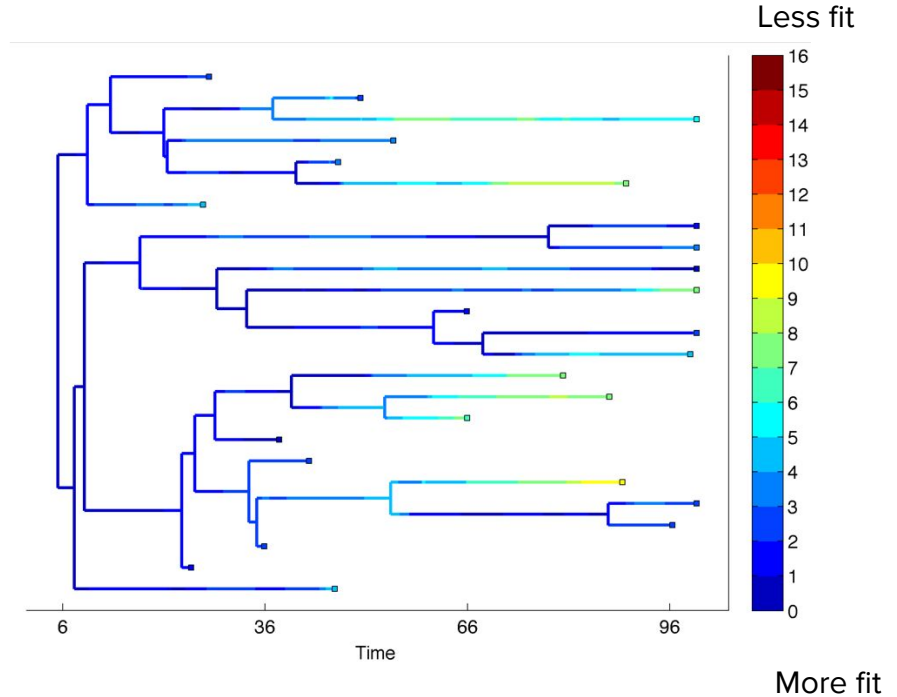
MTBD models allow us to compute the **joint likelihood** that both the tree and the observed tip genotypes evolved exactly as observed (Stadler and Bonhoeffer, 2013).



Fitness shapes trees

More fit lineages will be transmitted (branch) more often and leave behind more sampled descendants than less fit lineages.

Estimating transmission rates from the branching structure of phylogenies using MTBD provides us with one way to directly estimate pathogen fitness from genomic data.



Pathogen transmission fitness at the host-population level is determined by many different factors...

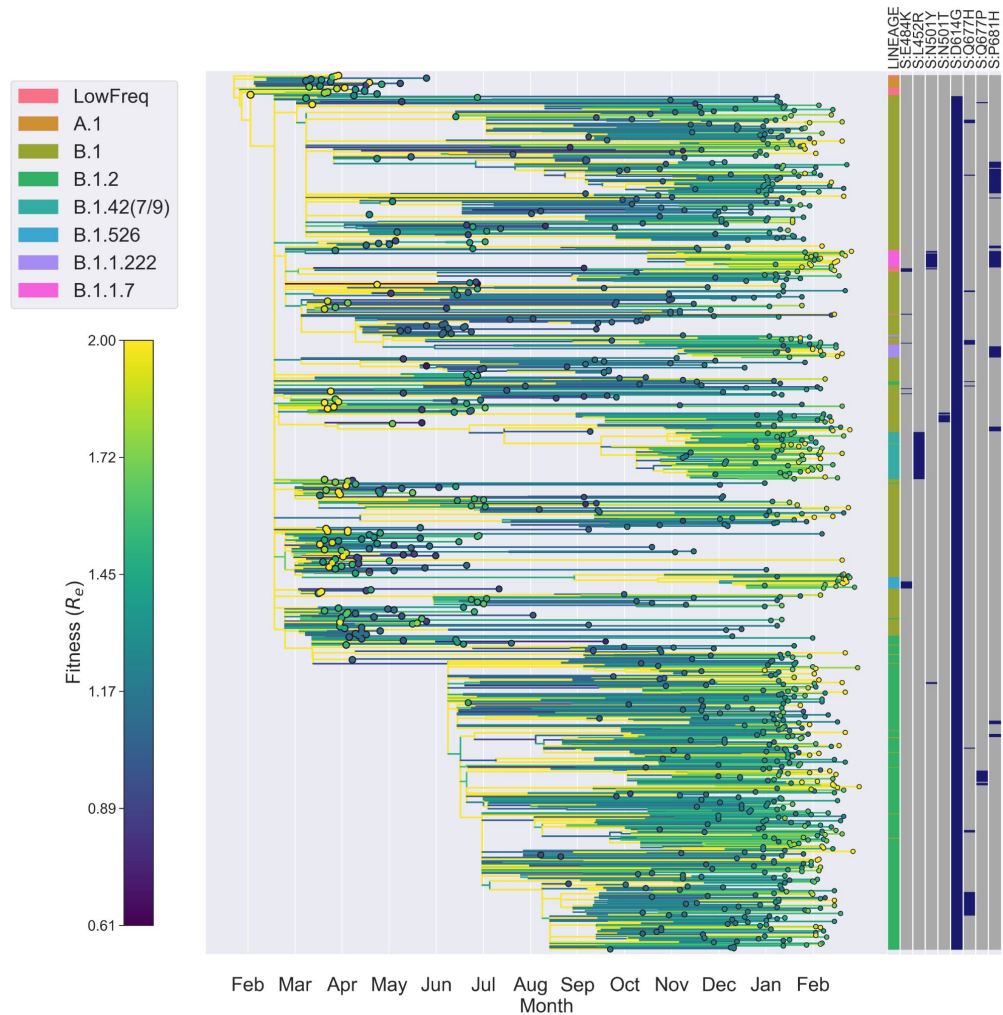
Can we use genomic data to learn what factors shape the transmission fitness of SARS-CoV-2?



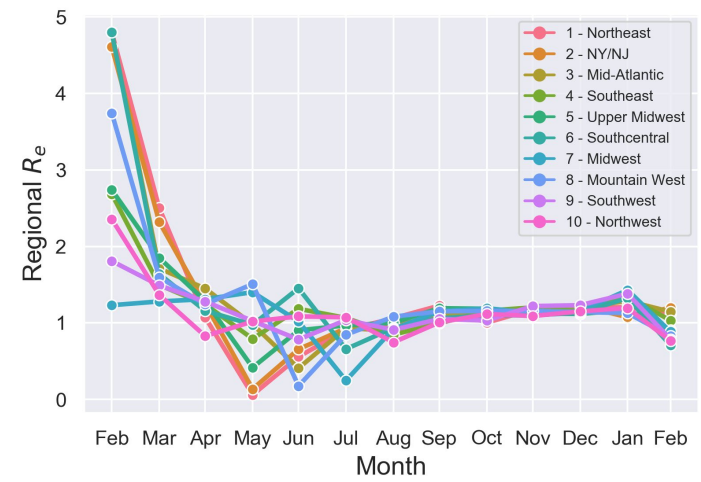
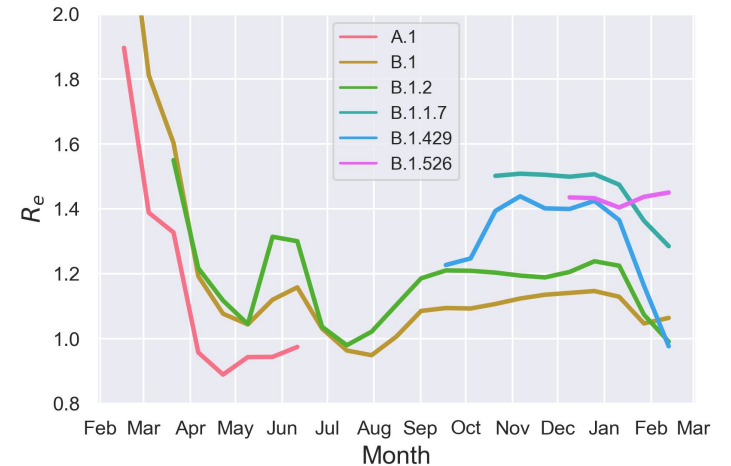
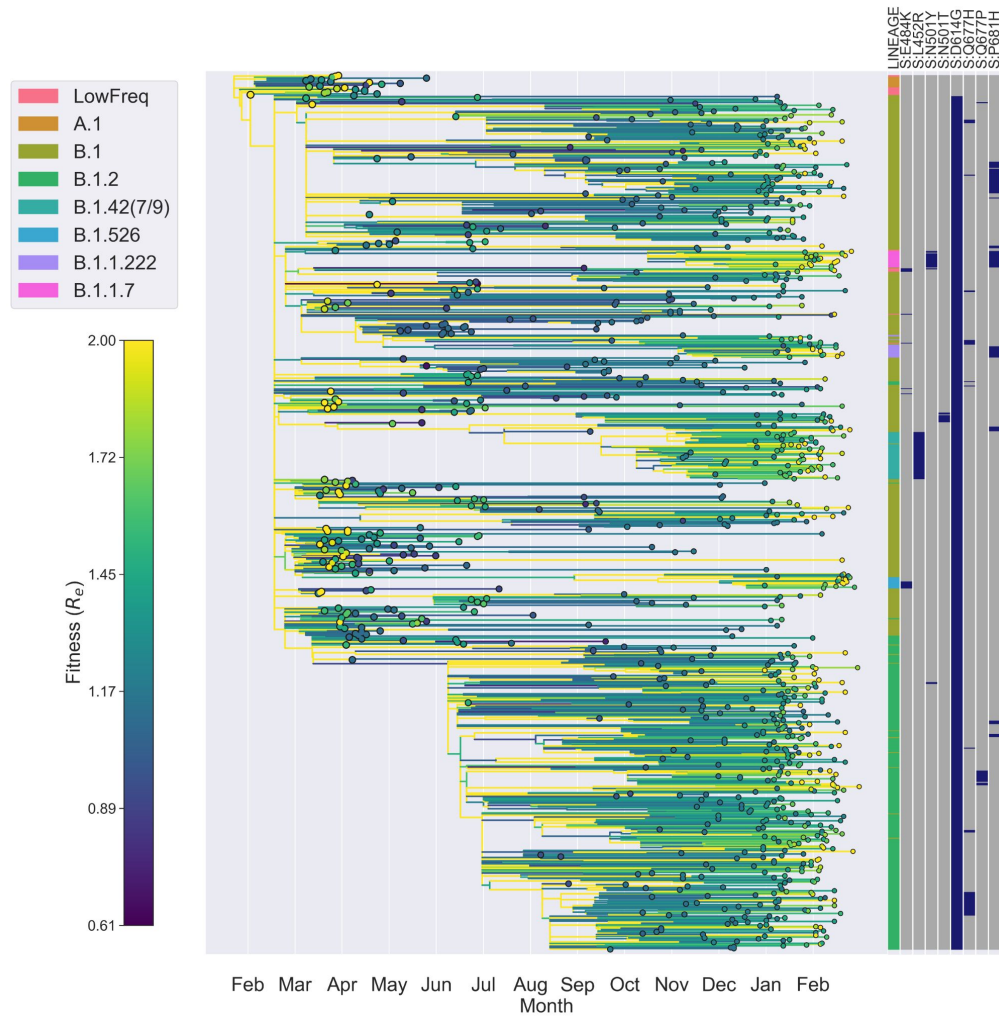
Lenora Kepler
(NCSU, Bioinformatics)



Marco Hamins-Puertolas
(NCSU, Biomathematics)



Kepler *et al.* (Virus Evolution, 2021)



SARS-CoV-2 workflow

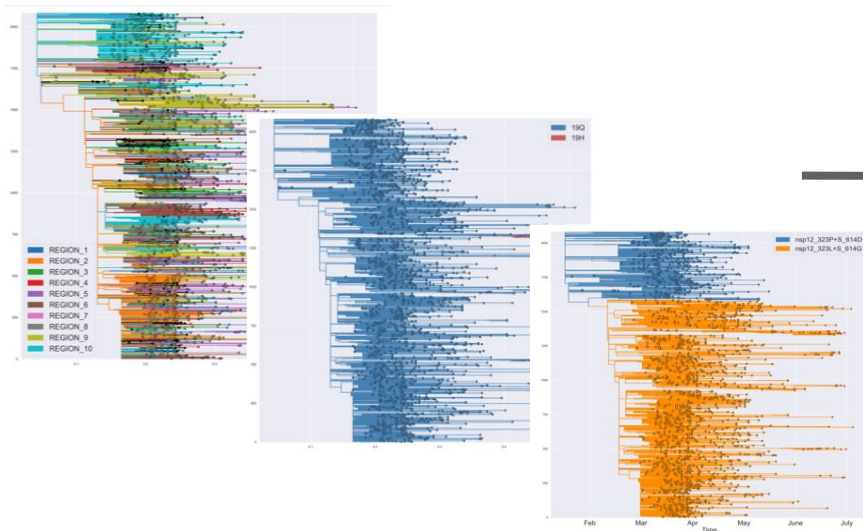
Sequences
download from
GISAID (n = 88,000)

Phylogenetic
reconstruction
(RAxML)

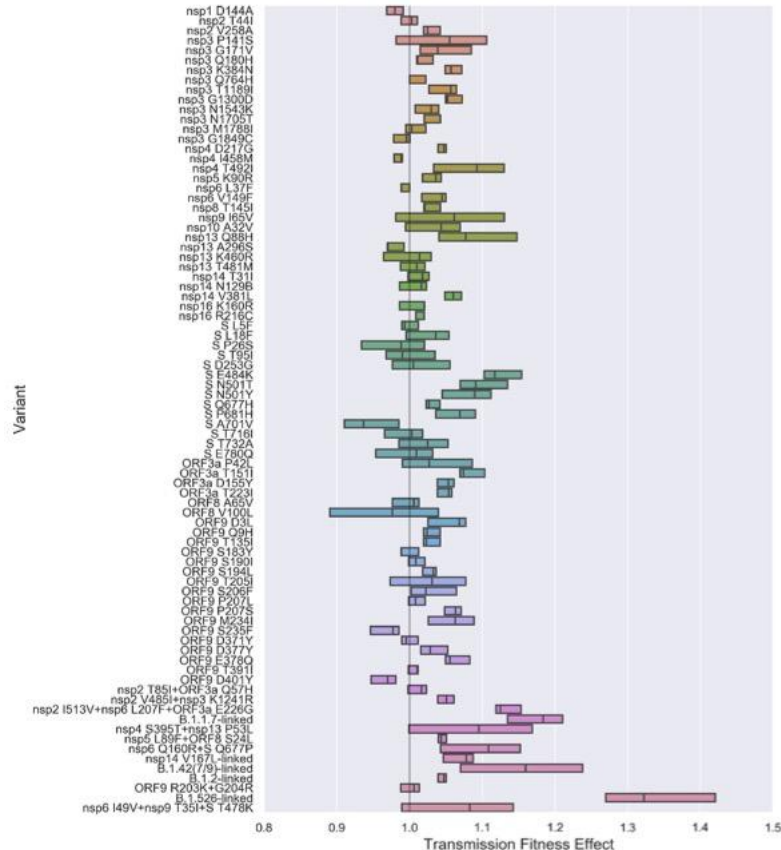
Least squares dating (LSD)

Ancestral feature
reconstruction
(PastML)

Features encoded as binary predictor variables

[illegible]

Fitness effects of amino acid variants



Variant	MLE	95% CI	Frequency
nsf3 K384N	1.057	1.05-1.07	0.016
nsf3 T1189I	1.057	1.03-1.06	0.012
nsf3 G1300D	1.052	1.05-1.07	0.023
nsf4 T429I	1.093	1.04-1.13	0.019
nsf13 Q88H	1.077	1.04-1.14	0.012
nsf14 V381L	1.06	1.05-1.07	0.014
ORF3a T151I	1.074	1.07-1.11	0.015
ORF3a D155Y	1.053	1.04-1.06	0.018
ORF3a T223I	1.053	1.04-1.06	0.033
ORF3a E226G	1.124	1.12-1.15	0.014
ORF9 D3L	1.068	1.03-1.08	0.015
ORF9 P207S	1.063	1.05-1.07	0.015
ORF9 M234I	1.063	1.03-1.086	0.054
ORF9 E378Q	1.055	1.05-1.08	0.012

Phylogenies can tell us about:

- Linkage and the sources of transmission
- The origins of epidemics and new strains
- Past epidemic dynamics
- Pathogen fitness and adaptation

**What do you want to
learn from this class?**

For Wednesday

On Wednesday we'll start with a tutorial that should help us ease into working with sequence data and trees.

Please have your laptops ready!

Try to install RAxML ahead of time

If you're interested in doing the Python exercises, install Python (with Anaconda) and Biopython.