

Introduction to phylogenetics

Louise Moncla
Genomics of Disease in Wildlife Workshop
June 8, 2023



Thijs Kuiken @thijskuiken · Feb 23

...

Following the death of an 11-year-old girl from highly pathogenic avian #influenza #H5N1 in Prey Veng province, #Cambodia, Ministry of Health reports 12 more infected people, 4 of whom have begun to show symptoms. Results of diagnosis expected tomorrow.



khmertimeskh.com

After death of girl, 12 more possibly detected with H5N1 bird flu in C...

Ms. Youk Sambath, Secretary of State of the Ministry of Health, has confirmed that the Ministry of Health's emergency response team h...



55



1,023

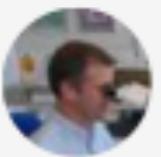


1,507



625K

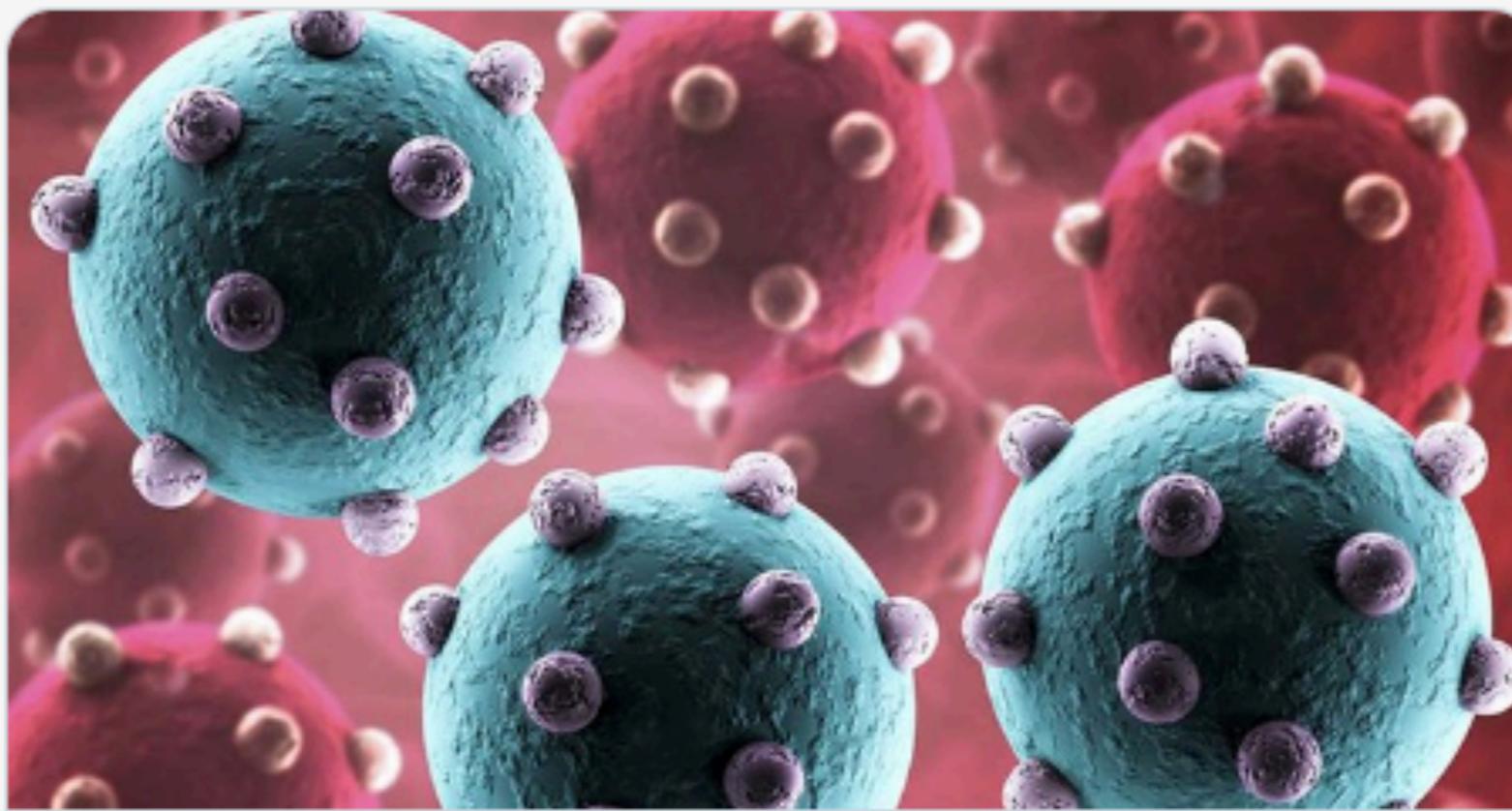




Thijs Kuiken @thijskuiken · Feb 23

...

Following the death of an 11-year-old girl from highly pathogenic avian #influenza #H5N1 in Prey Veng province, #Cambodia, Ministry of Health reports 12 more infected people, 4 of whom have begun to show symptoms. Results of diagnosis expected tomorrow.



khmertimeskh.com

After death of girl, 12 more possibly detected with H5N1 bird flu in C...

Ms. Youk Sambath, Secretary of State of the Ministry of Health, has confirmed that the Ministry of Health's emergency response team h...

55

1,023

1,507

625K



Jurre Y Siegers, PhD @jurreysi · Feb 23

...

There seems to be a translation error from Khmer to English. The report below (in khmer) states there are 12 contacts to the index case and 4 amongst those 12 report flu like symptoms.

The questions: Is this human case related to the ongoing outbreak in Europe/the Americas? Does this reflect a new epidemiological pattern?



Erik Karlsson
@E_A_Karlsson

...

Happy to announce that the full genome sequence of the Cambodian #H5N1 #avianflu case is now available on @GISAID. This has been an incredible effort by a number of people: (1/)

12:15 AM · Feb 26, 2023 · 532.9K Views

387 Retweets 74 Quotes 1,615 Likes 85 Bookmarks



Phylogeny



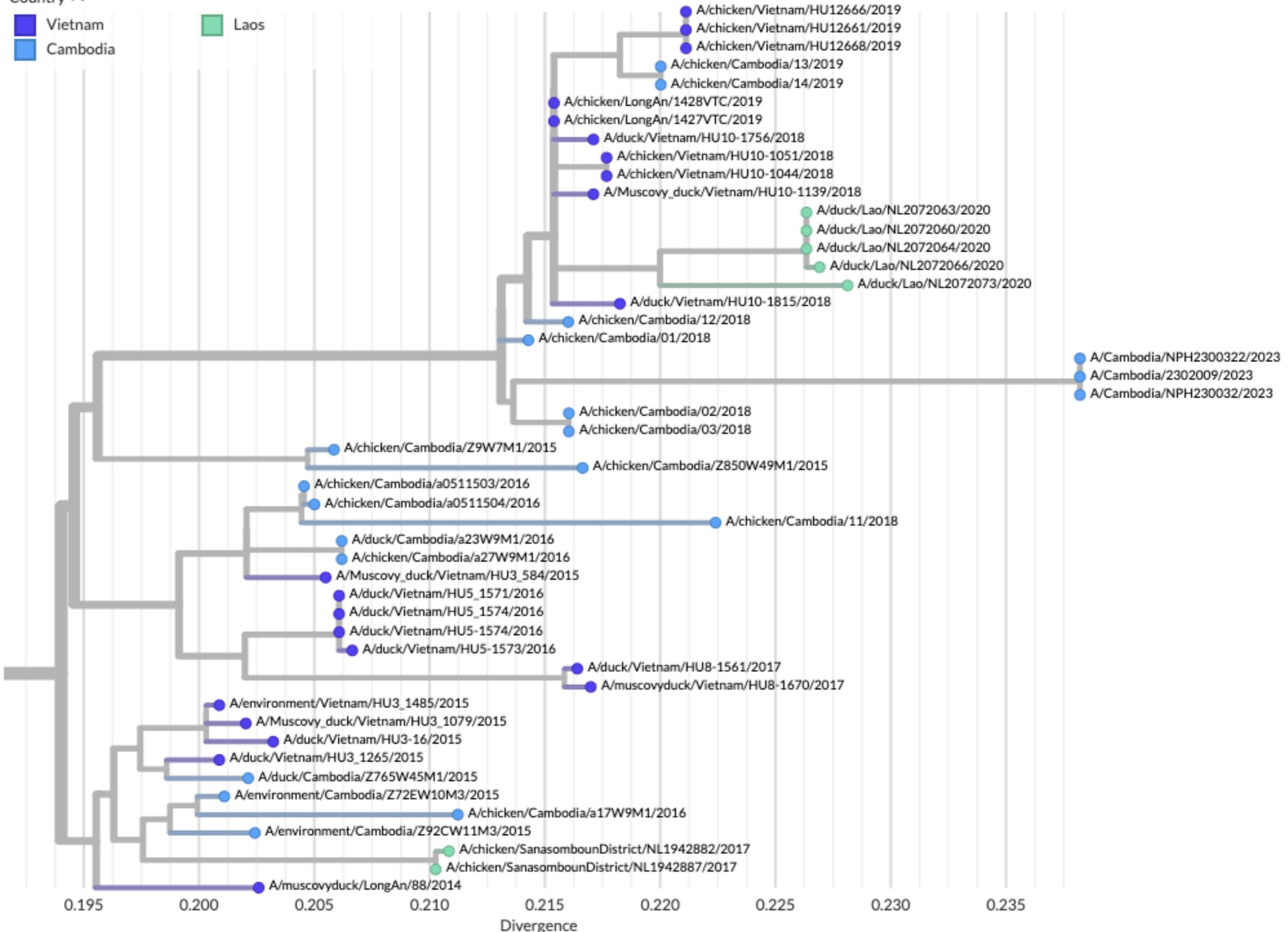
ZOOM TO SELECTED

RESET LAYOUT

Country ▲

- Vietnam
Cambodia

Laos



0.195 0.200 0.205 0.210 0.215 0.220 0.225 0.230 0.235

Divergence

Phylogeny

Subtype ^

h5n1
h5n6

h5n8

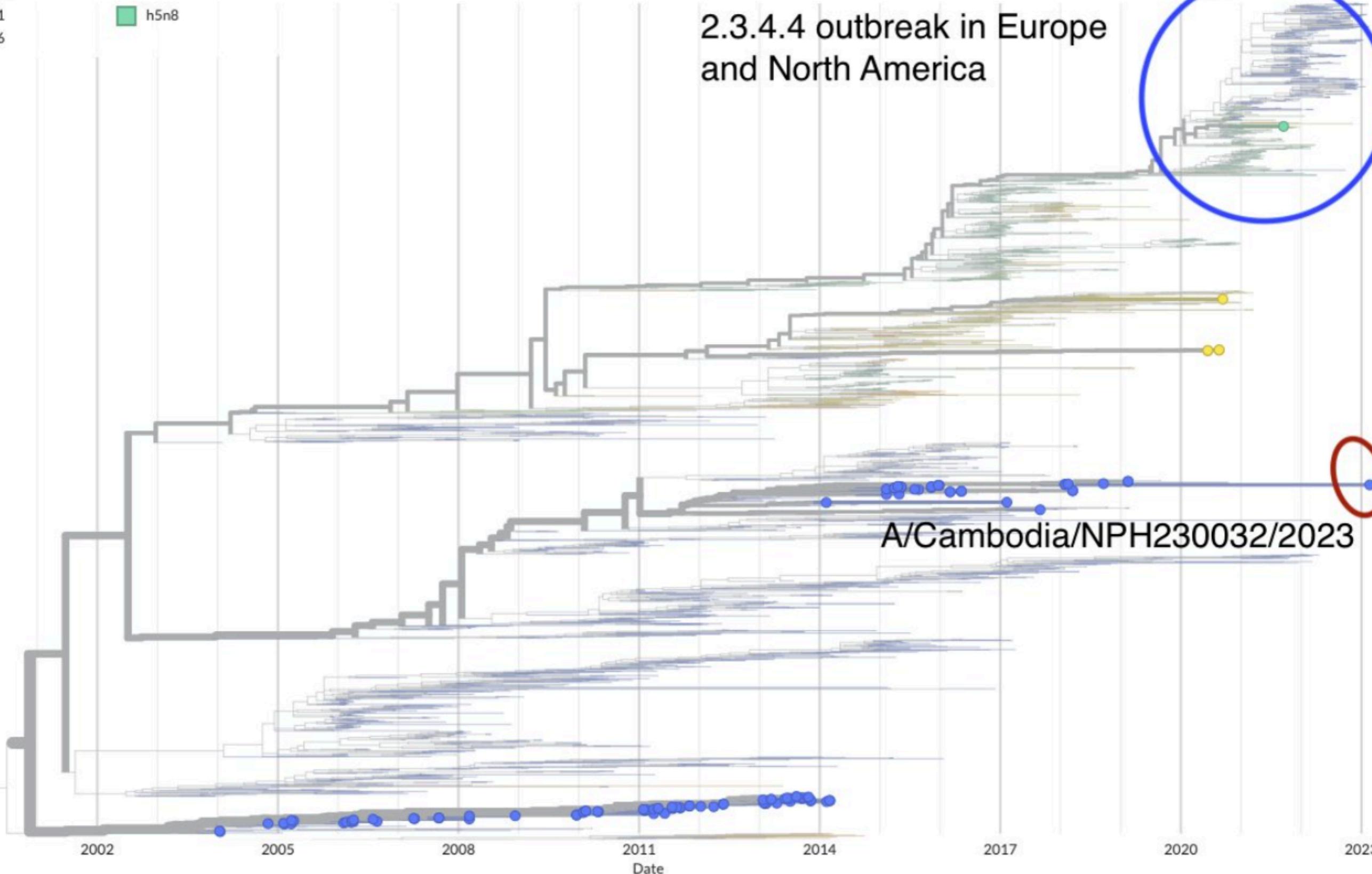


ZOOM TO SELECTED

RESET LAYOUT

2.3.4.4 outbreak in Europe and North America

A/Cambodia/NPH230032/2023



The answers:

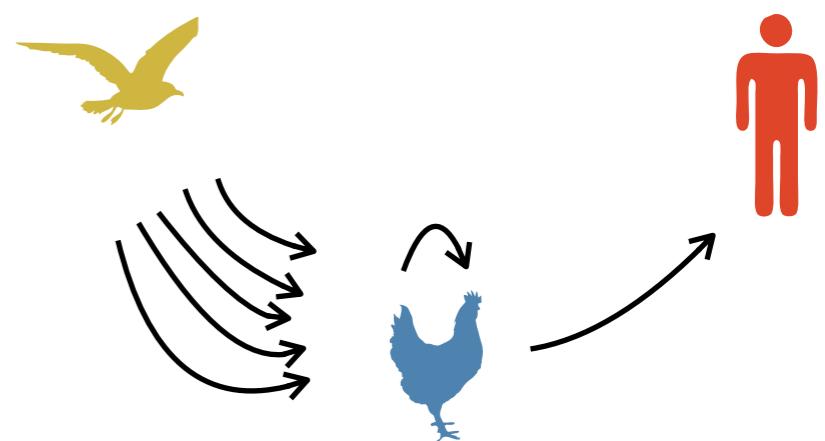
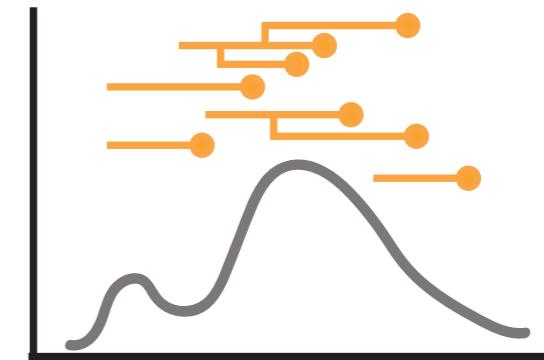
- These cases were caused by a lineage endemic to Cambodia, likely by direct interaction with sick, wild birds
- These cases do not represent evolution of human to human transmission of 2.3.4.4b viruses

Phylogenetics is the study of how epidemiological, immunological, and evolutionary processes act and interact to shape viral phylogenies.

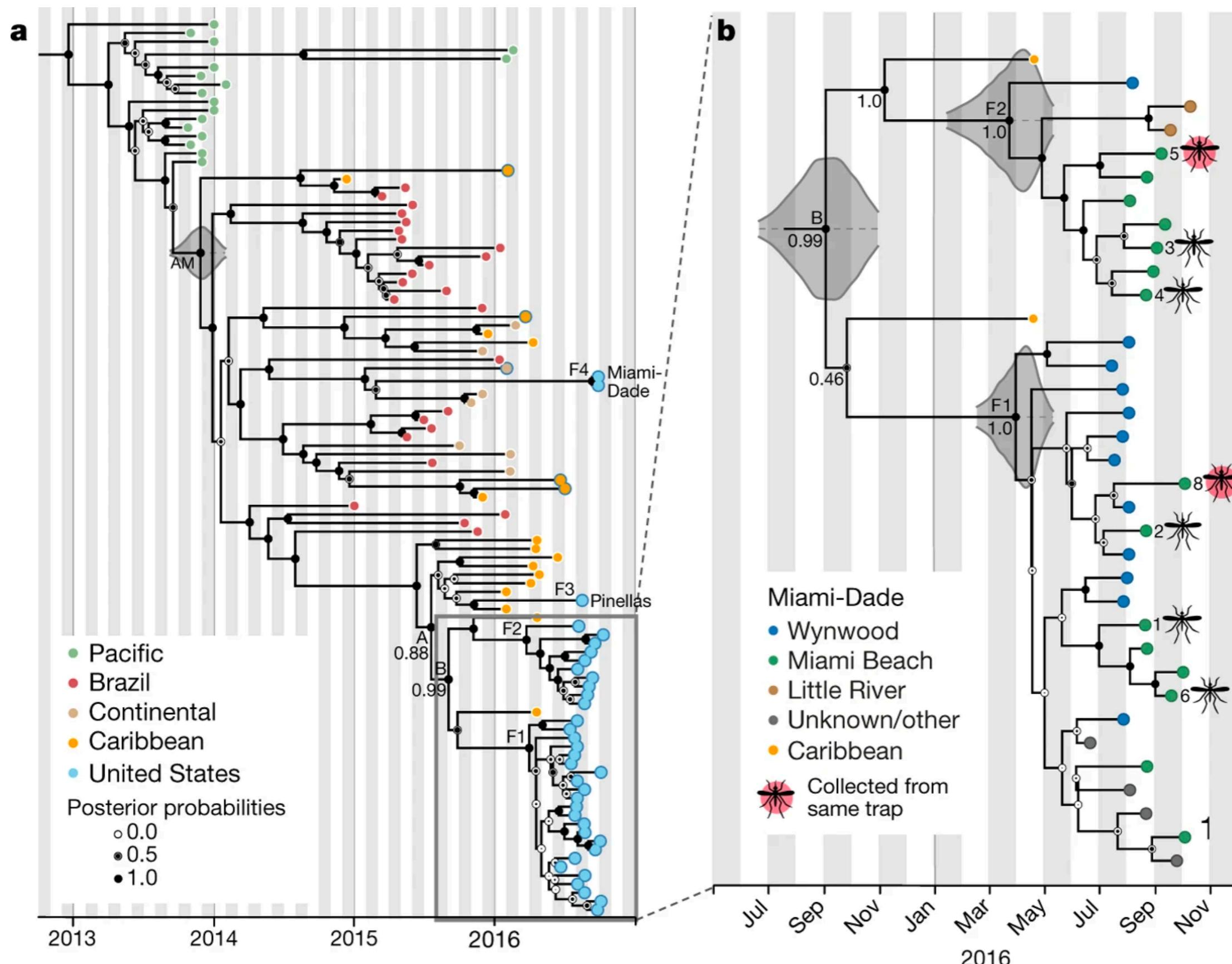
* Volz, Koelle, and Bedford, PLOS Computational Biology, 2013

What types of questions can phylodynamics answer?

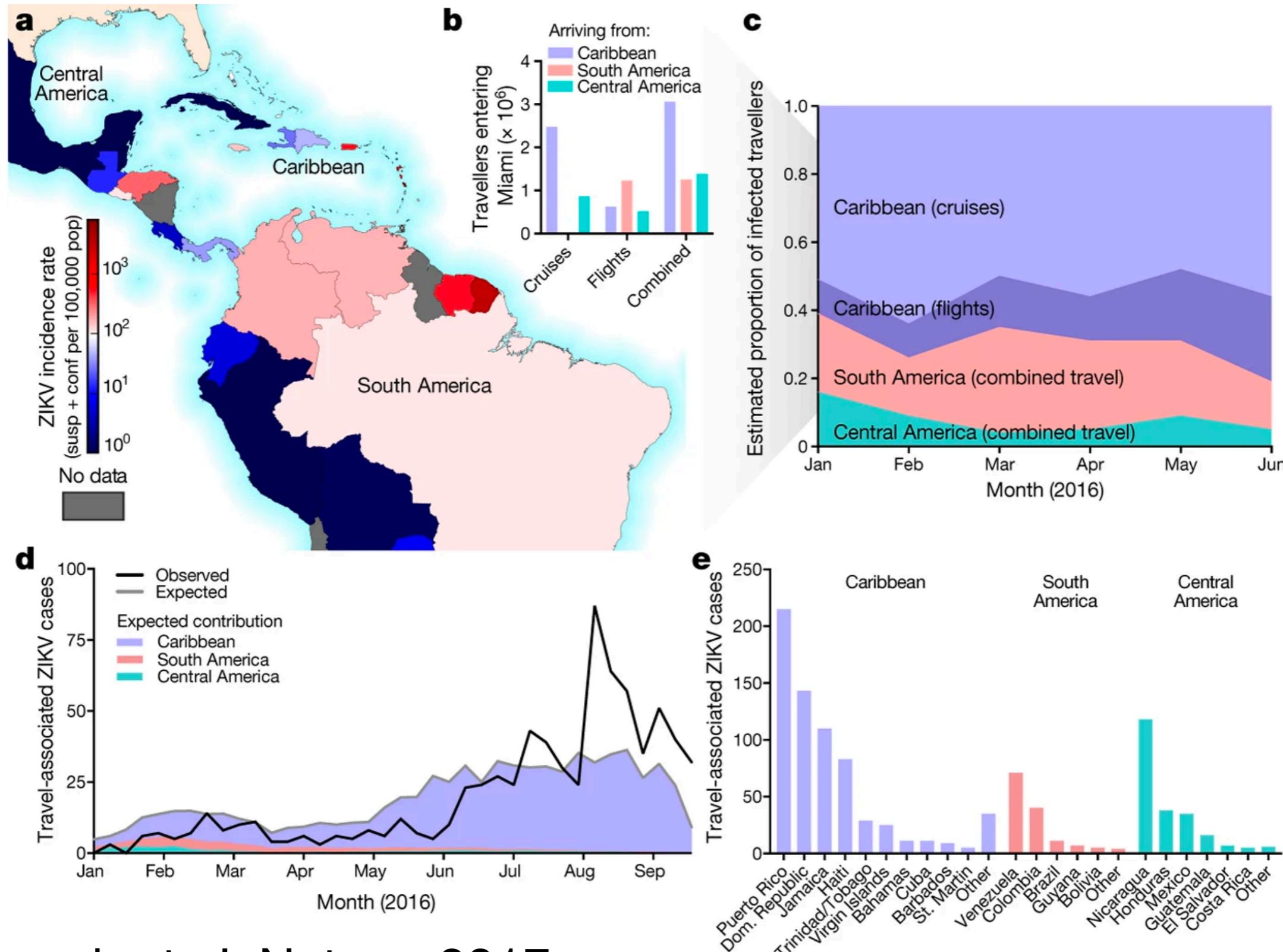
1. When did an outbreak begin, and where? How rapidly is it growing? Who is driving transmission?
2. Who infected whom?
3. How are pathogen populations shaped by natural selection?
4. How are pathogen populations moving across time, space, and hosts?



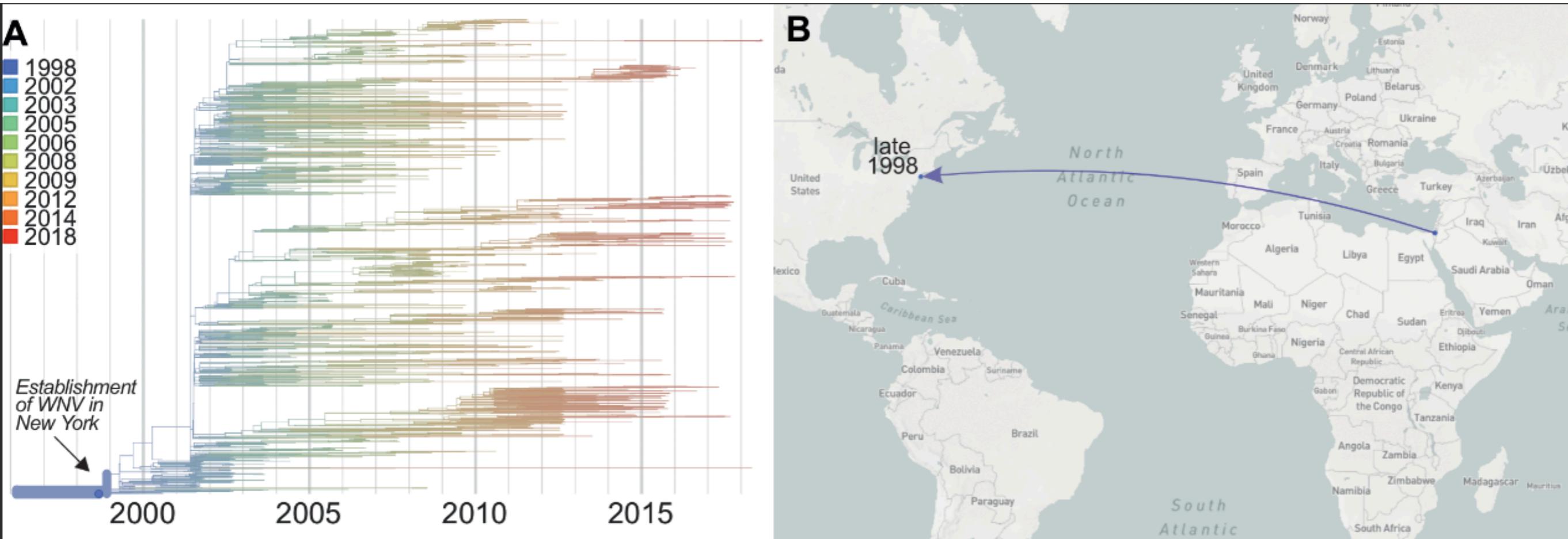
Inferring Zika introductions into Florida



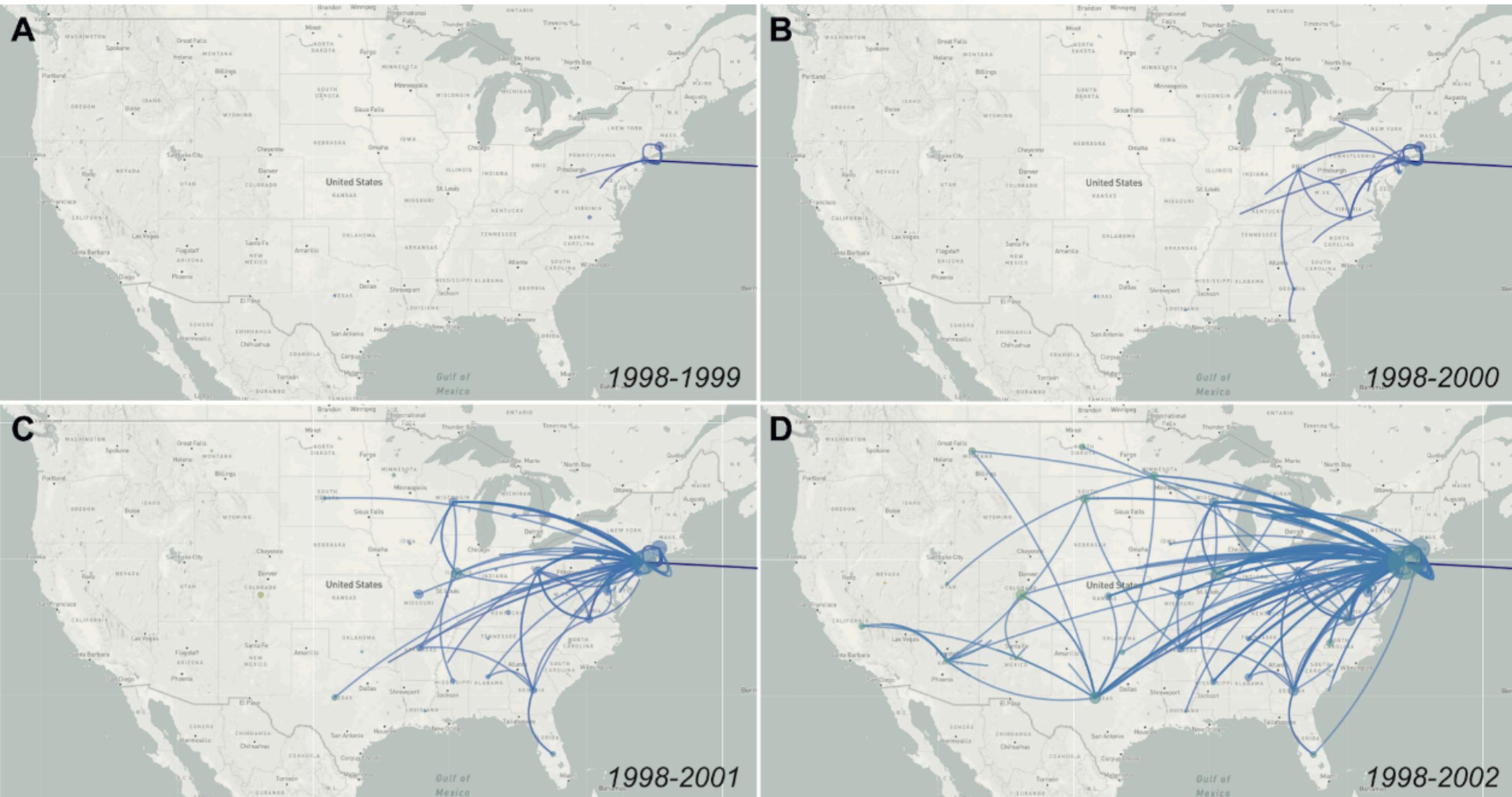
Mosquito abundance and cruise ship traffic made Miami uniquely at risk for Zika



West Nile virus was first introduced into North America in the late 1990s



West Nile virus spread from the east coast

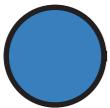


The key points we'll be making today:

1. Pathogen evolution happens on the same timescale as transmission, meaning that for rapidly evolving pathogens, **evolution, ecology, and epidemiology are linked**.
2. **Coalescent theory** allows us to generate expectations for phylogenies generated by populations that are evolving in particular ways.
3. These features allow us to use trees to infer parameters, and vice versa. We will see **examples** of how others have done this.
4. How can **we apply** these concepts to our own questions and data?

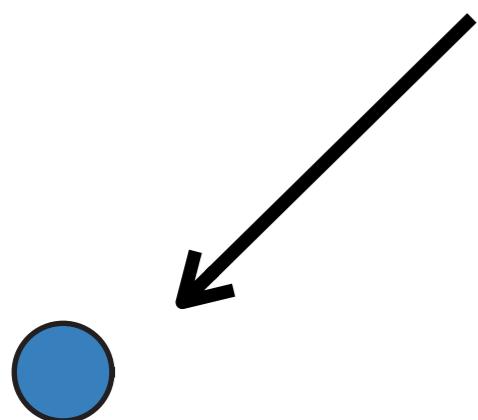
Phylogenetics is predicated on the fact that for rapidly evolving pathogens, ecological, evolutionary, and epidemiological dynamics all occur on the same timescale.

Viral genomes contain records of a virus's transmission history



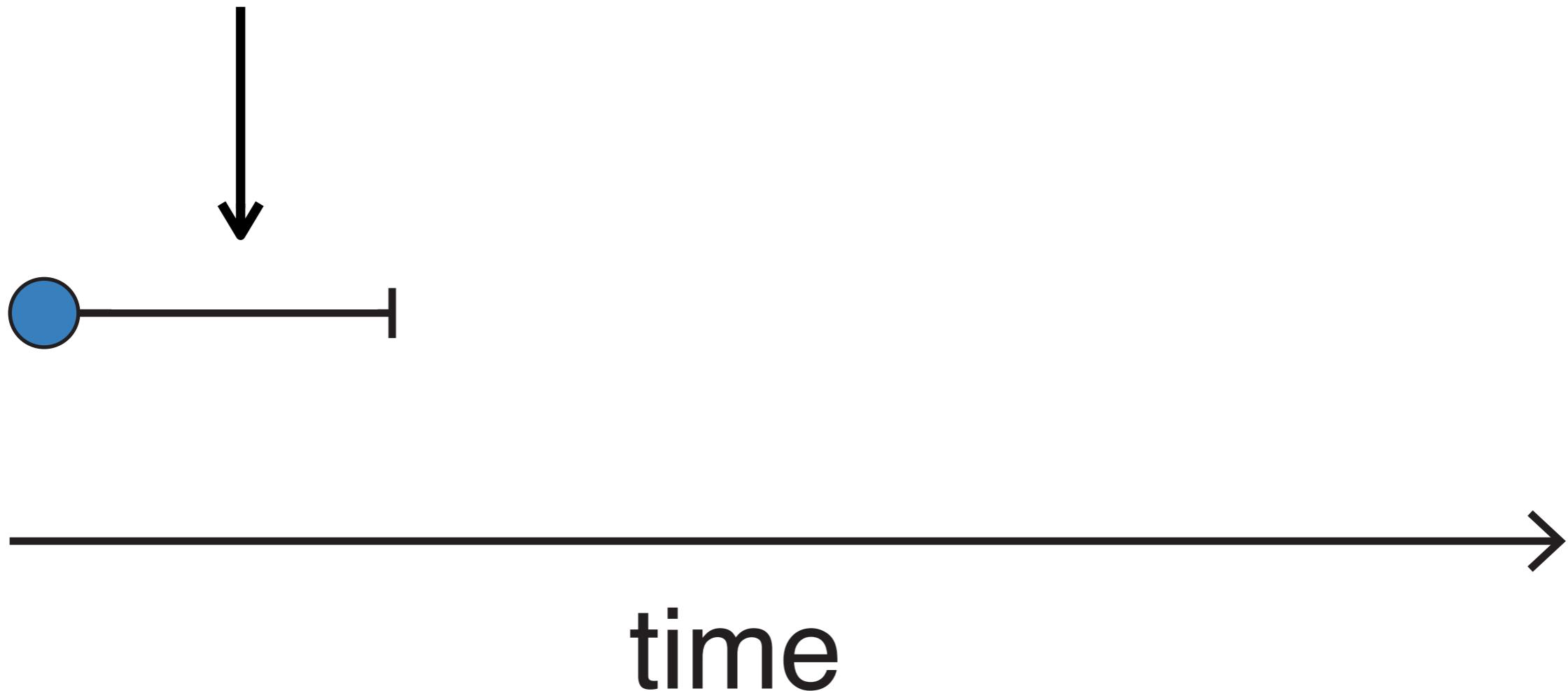
Viral genomes contain records of a virus's transmission history

one person
infected with
influenza

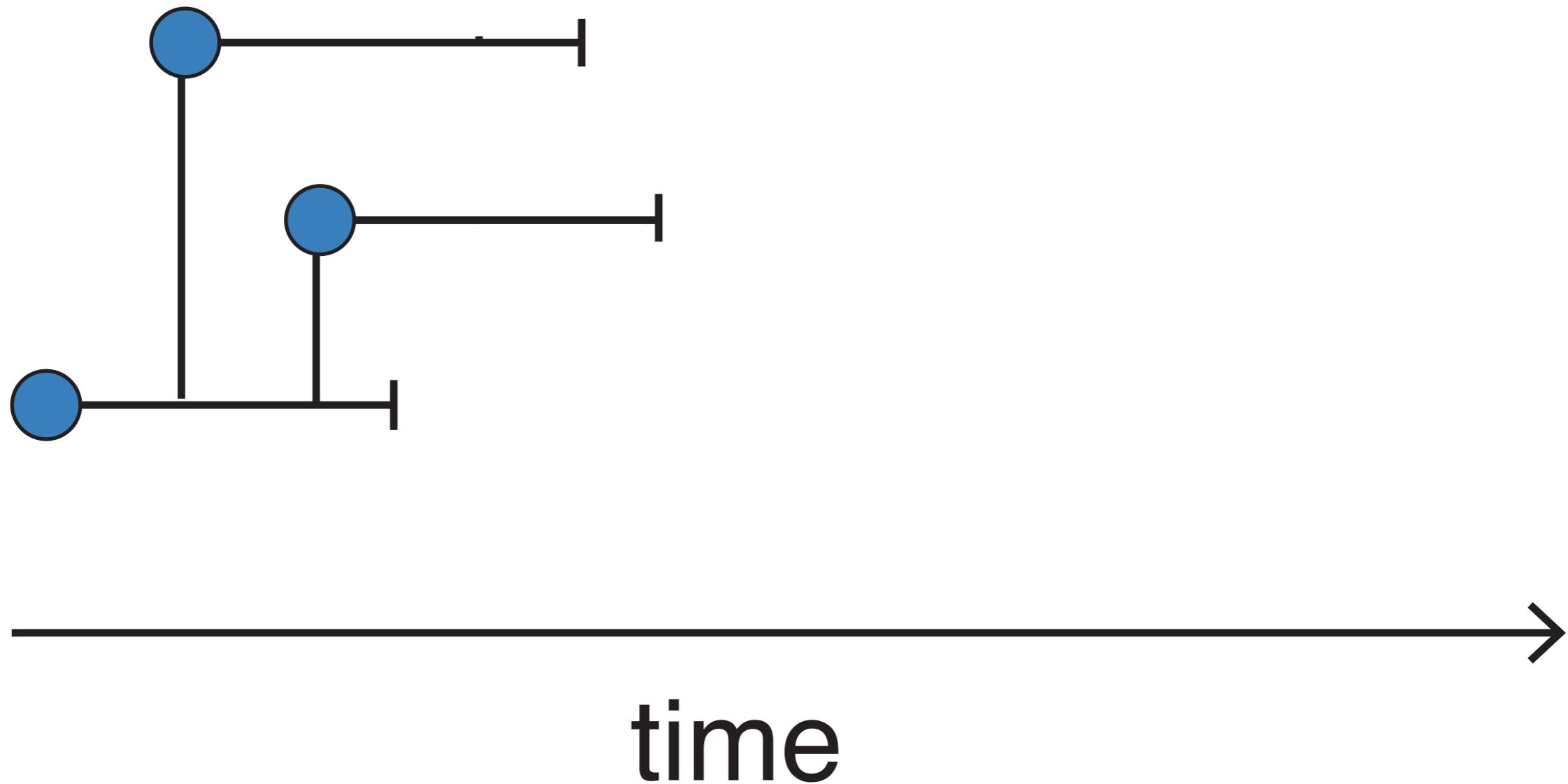


Viral genomes contain records of a virus's transmission history

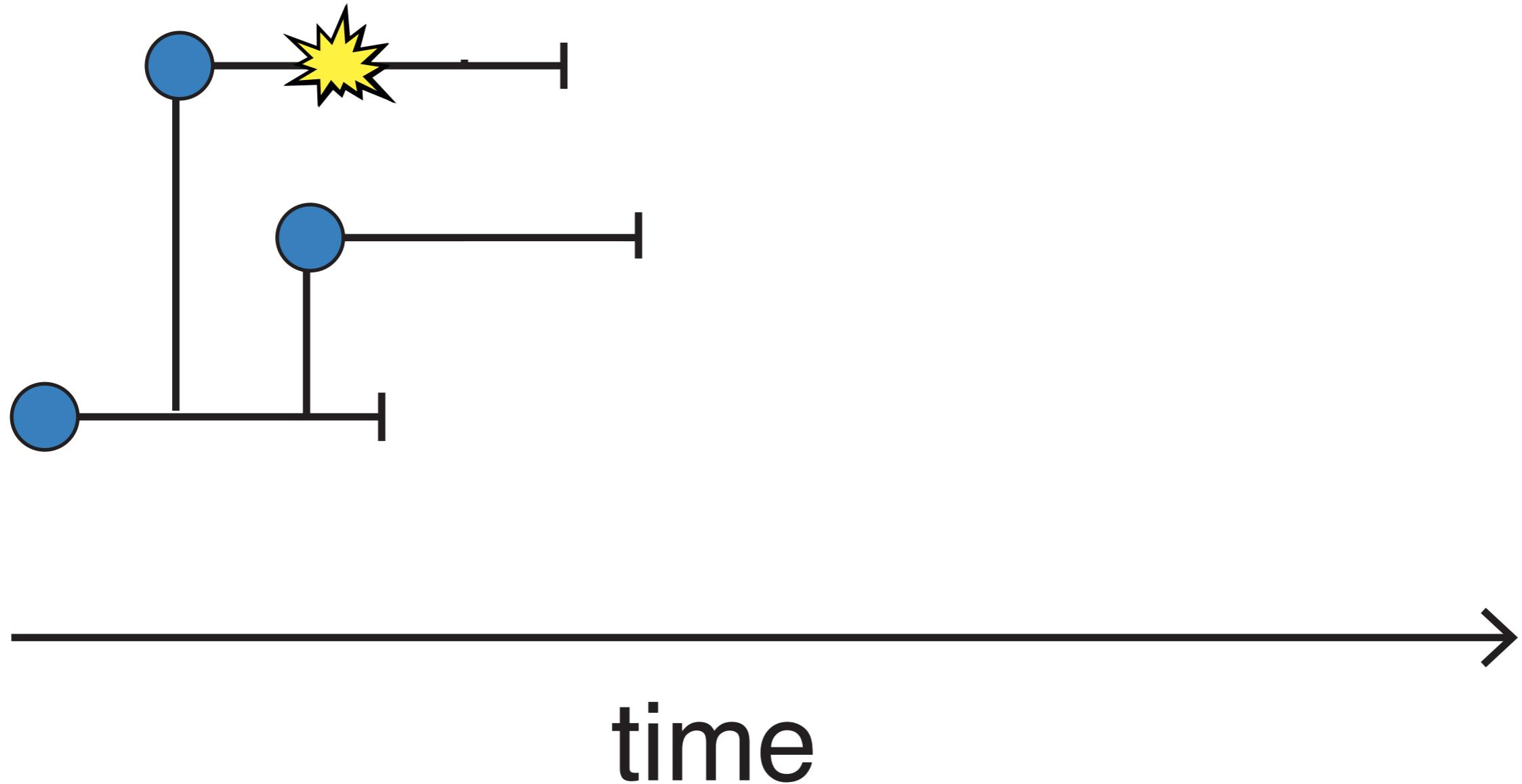
this person's
infectious period



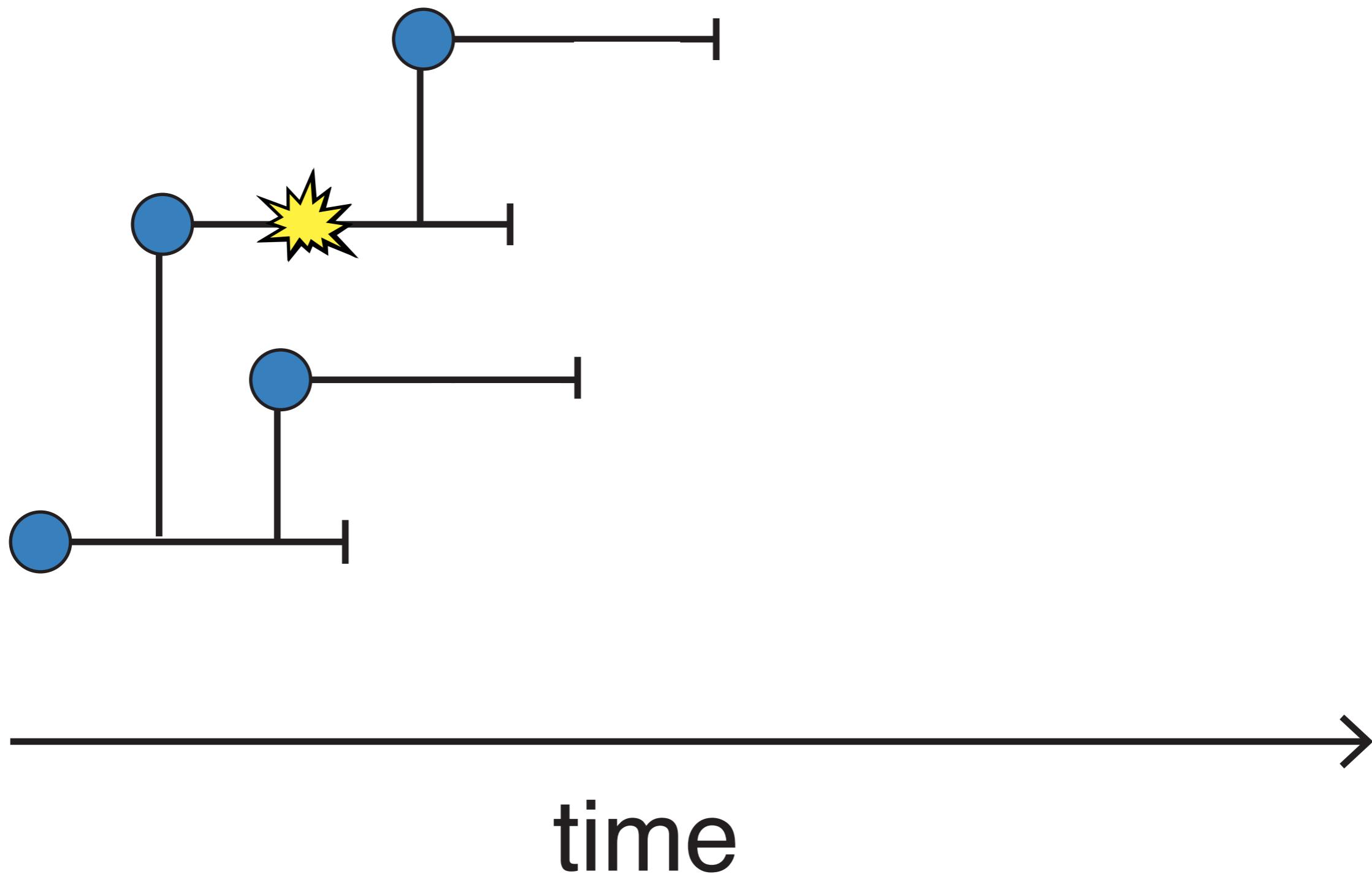
Viral genomes contain records of a virus's transmission history



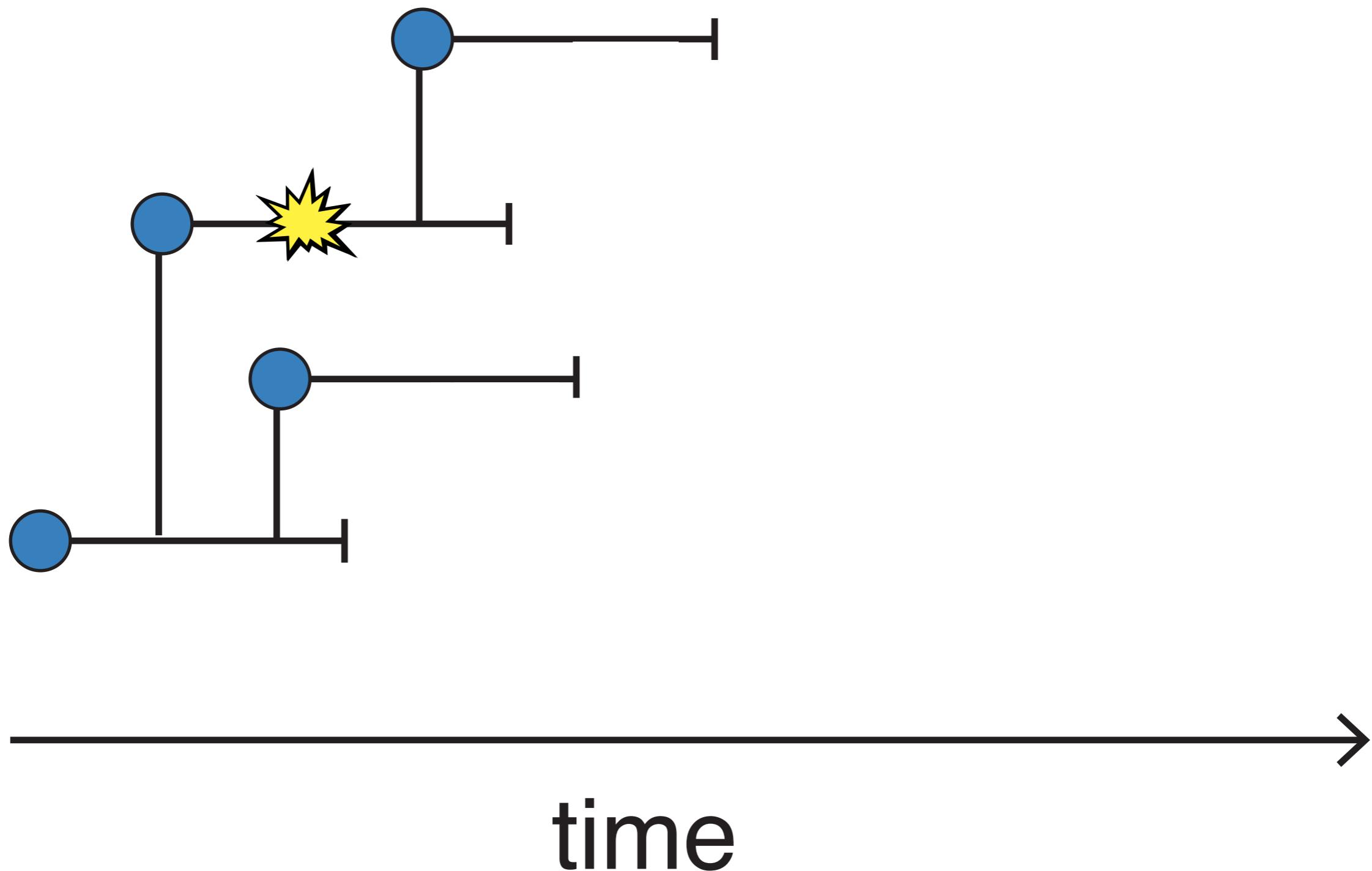
Viral genomes contain records of a virus's transmission history



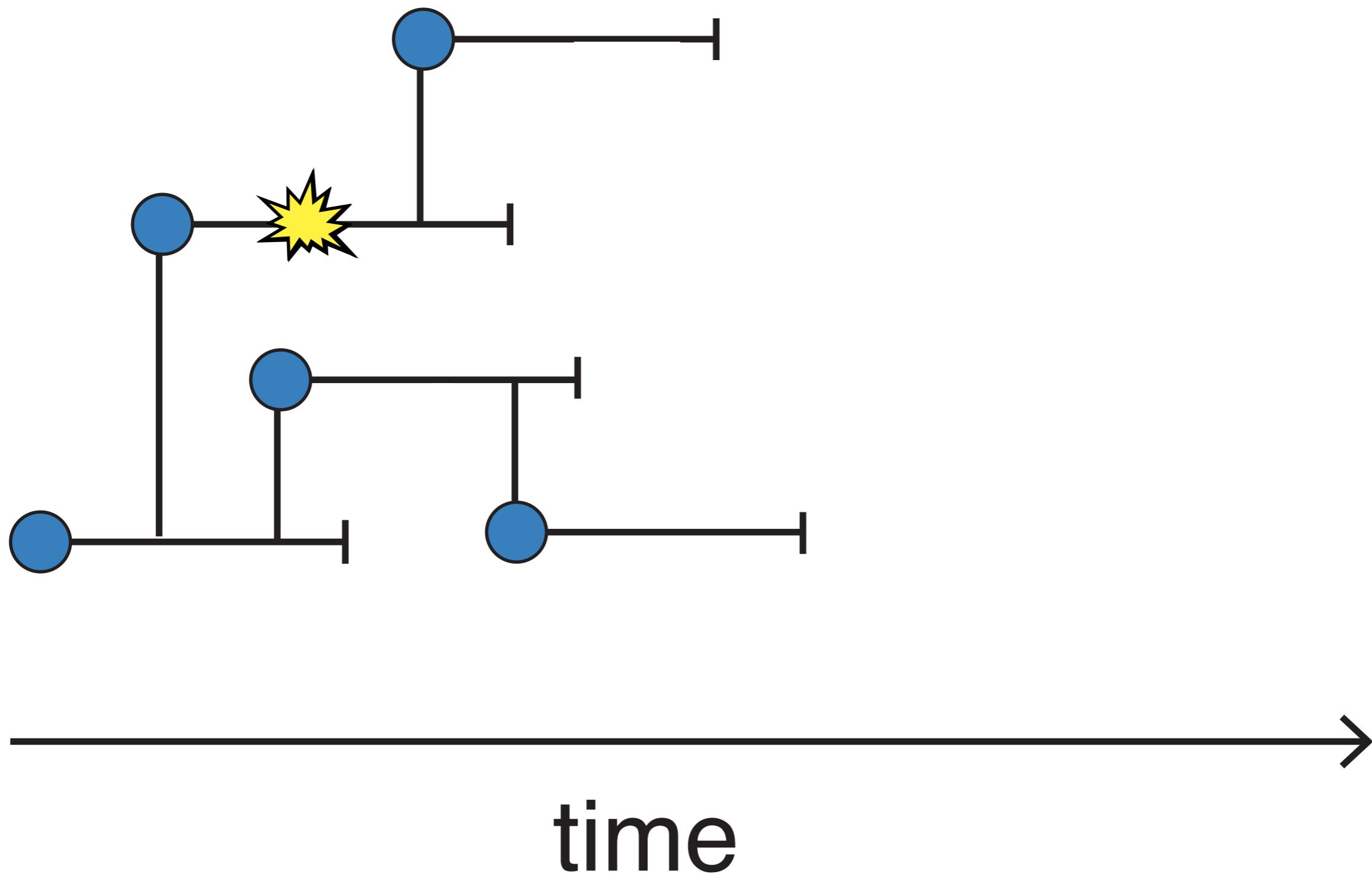
Viral genomes contain records of a virus's transmission history



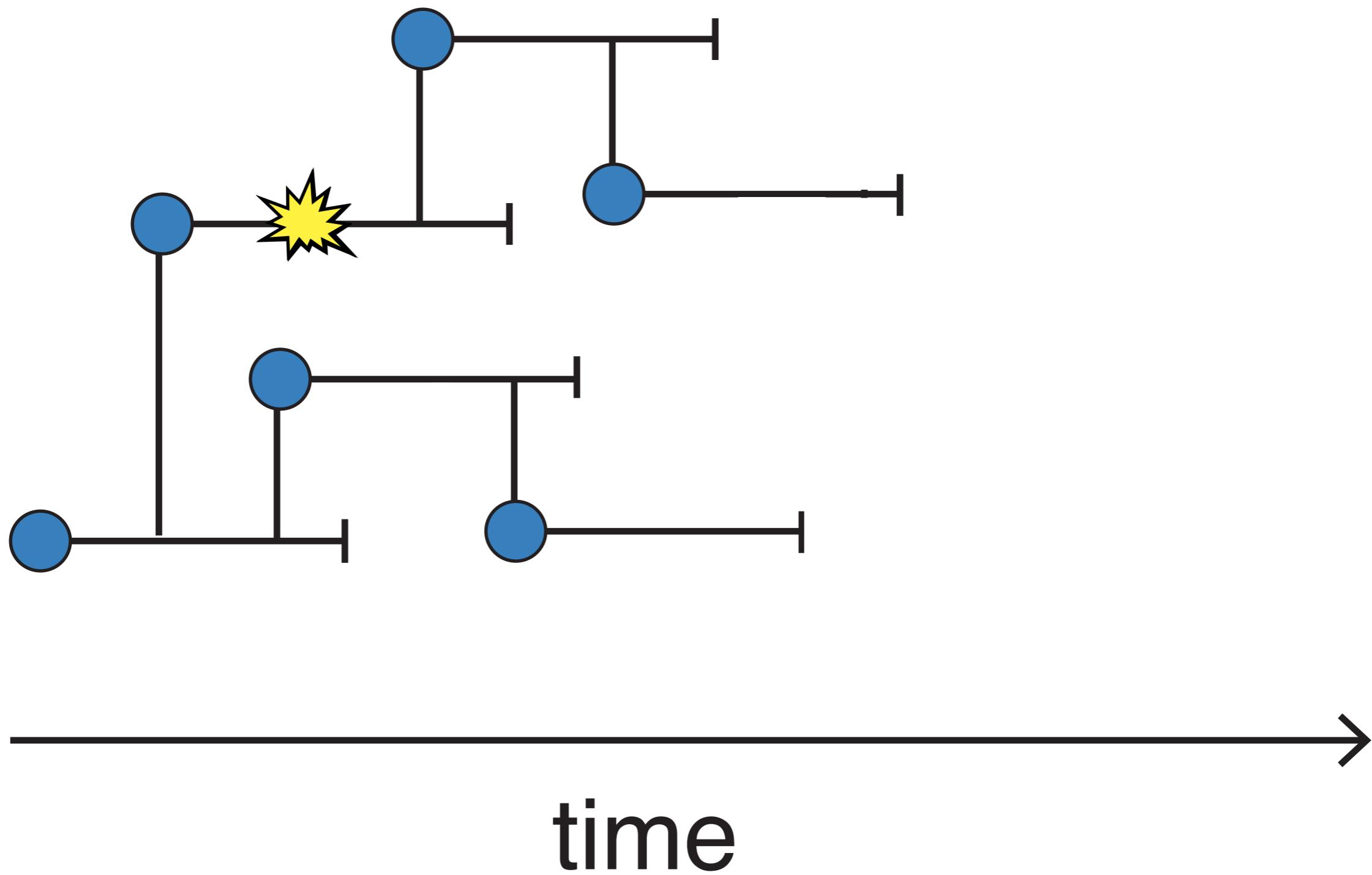
Viral genomes contain records of a virus's transmission history



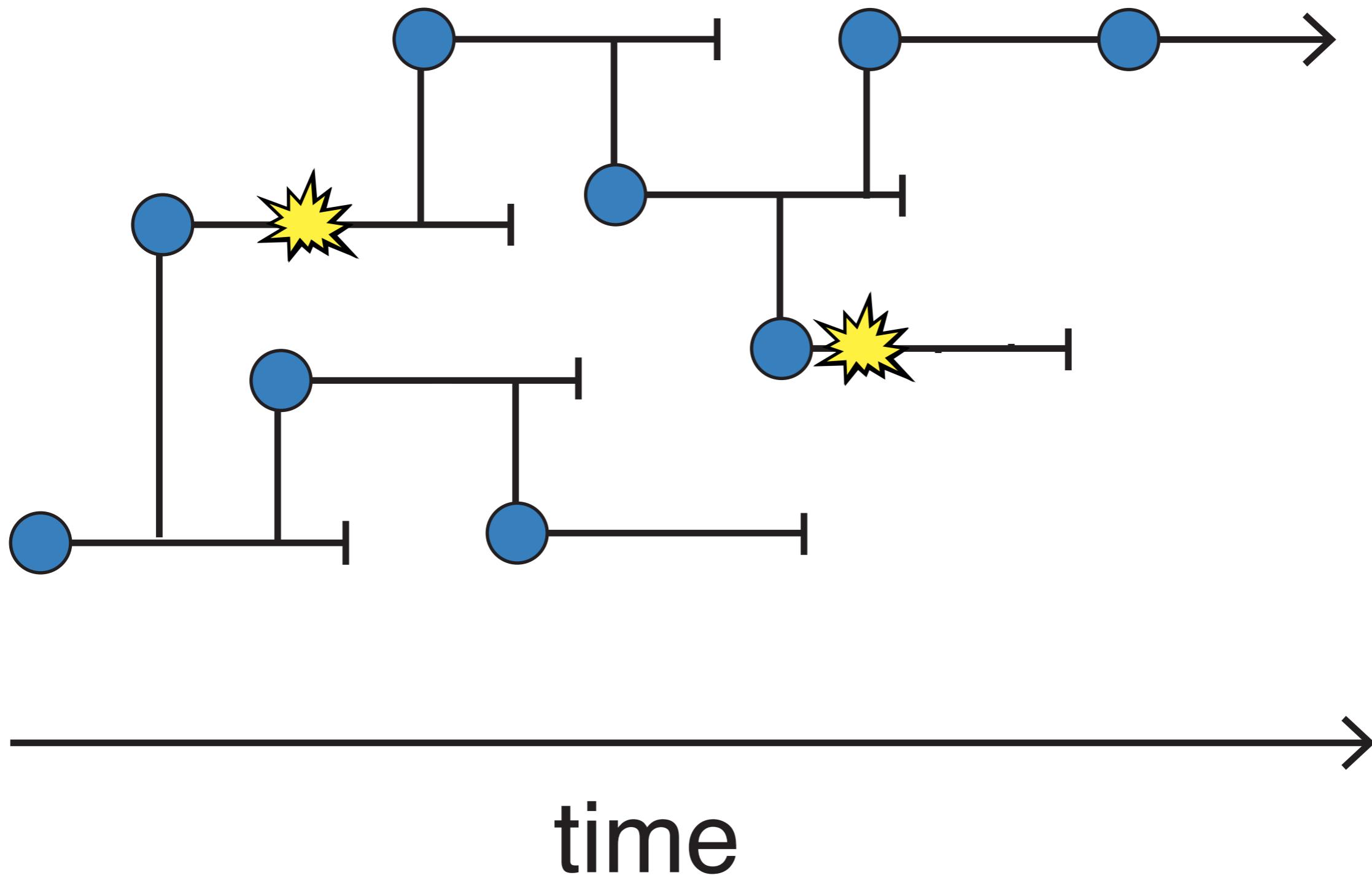
Viral genomes contain records of a virus's transmission history



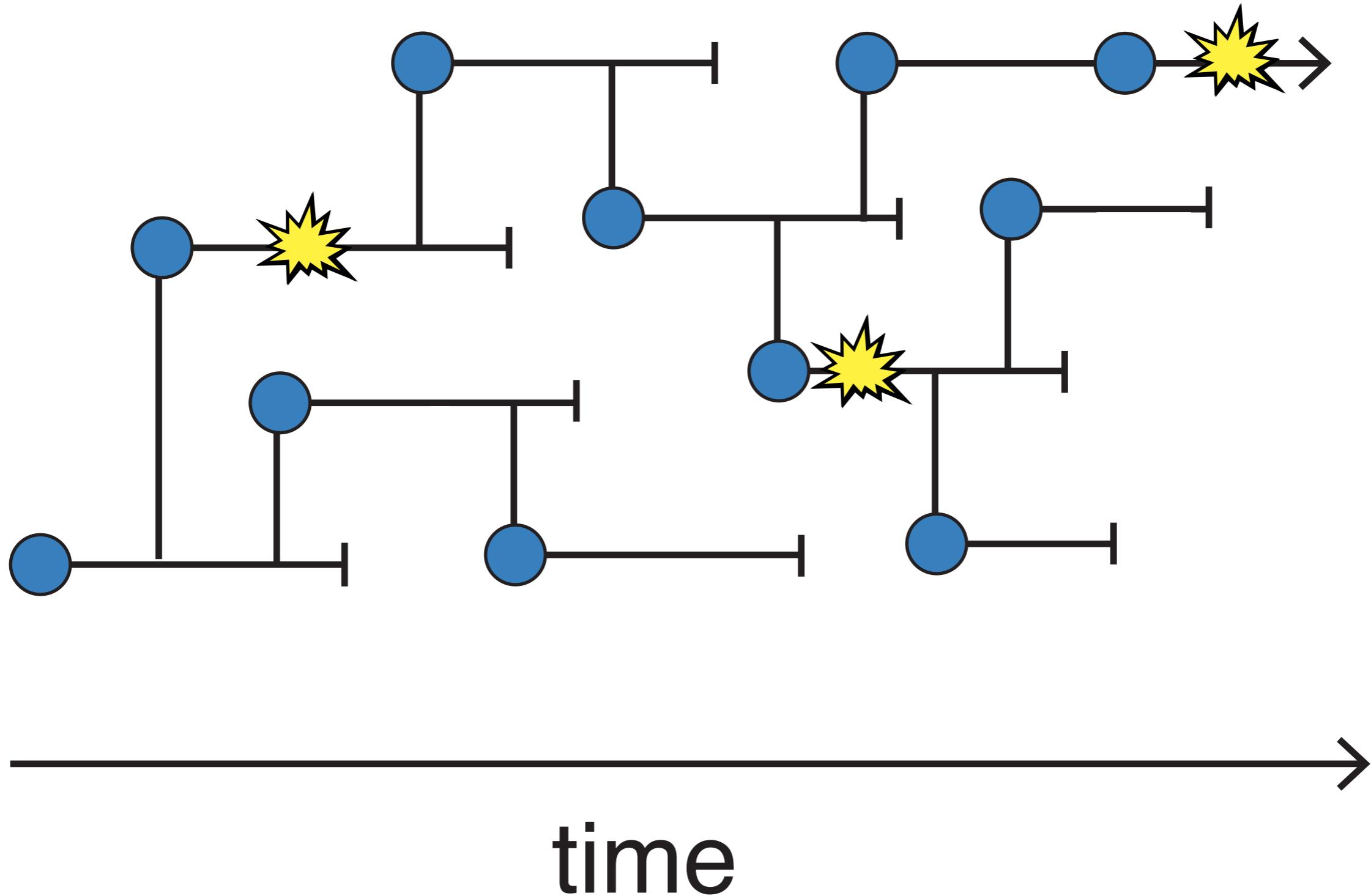
Viral genomes contain records of a virus's transmission history



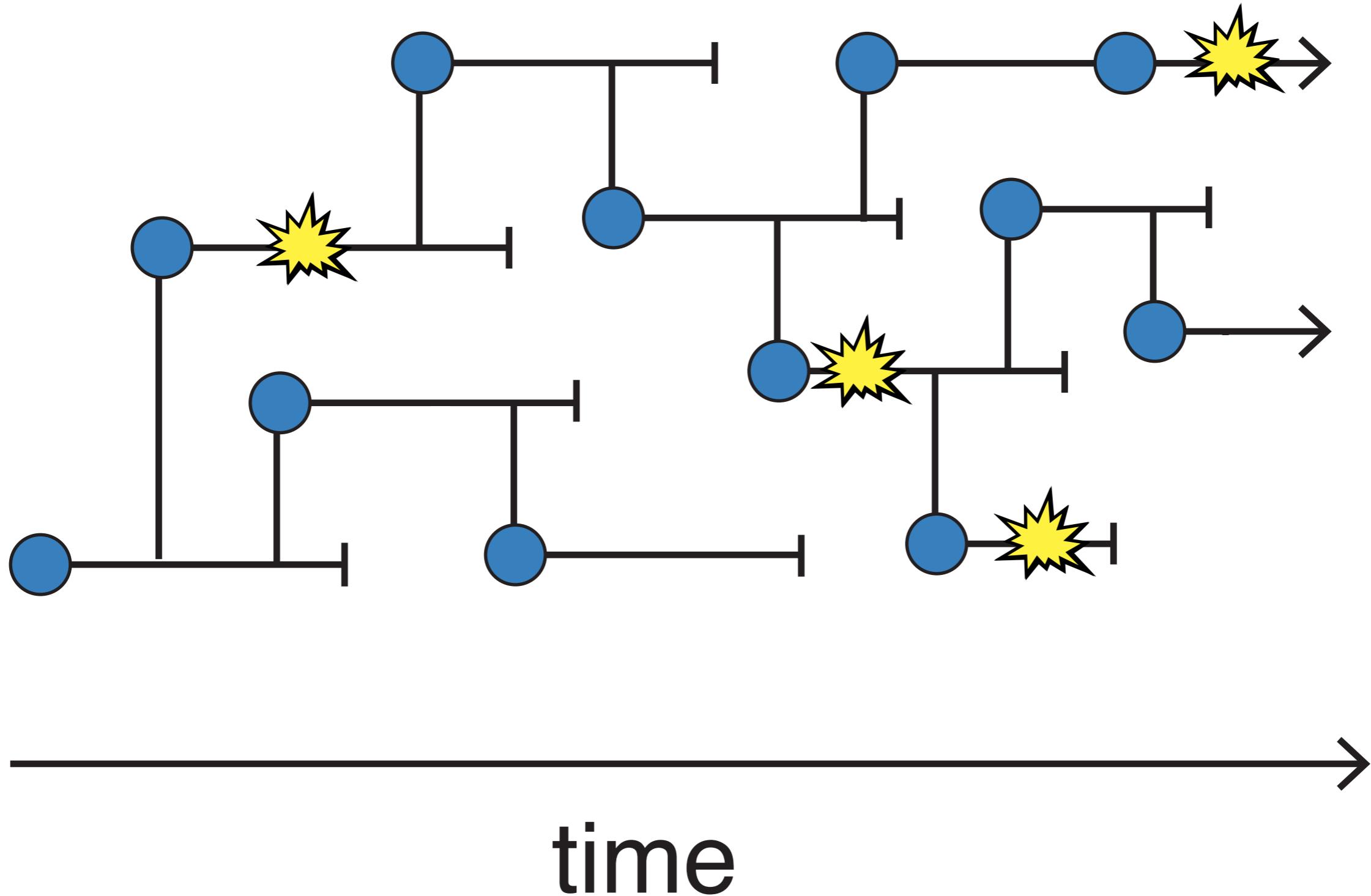
Viral genomes contain records of a virus's transmission history



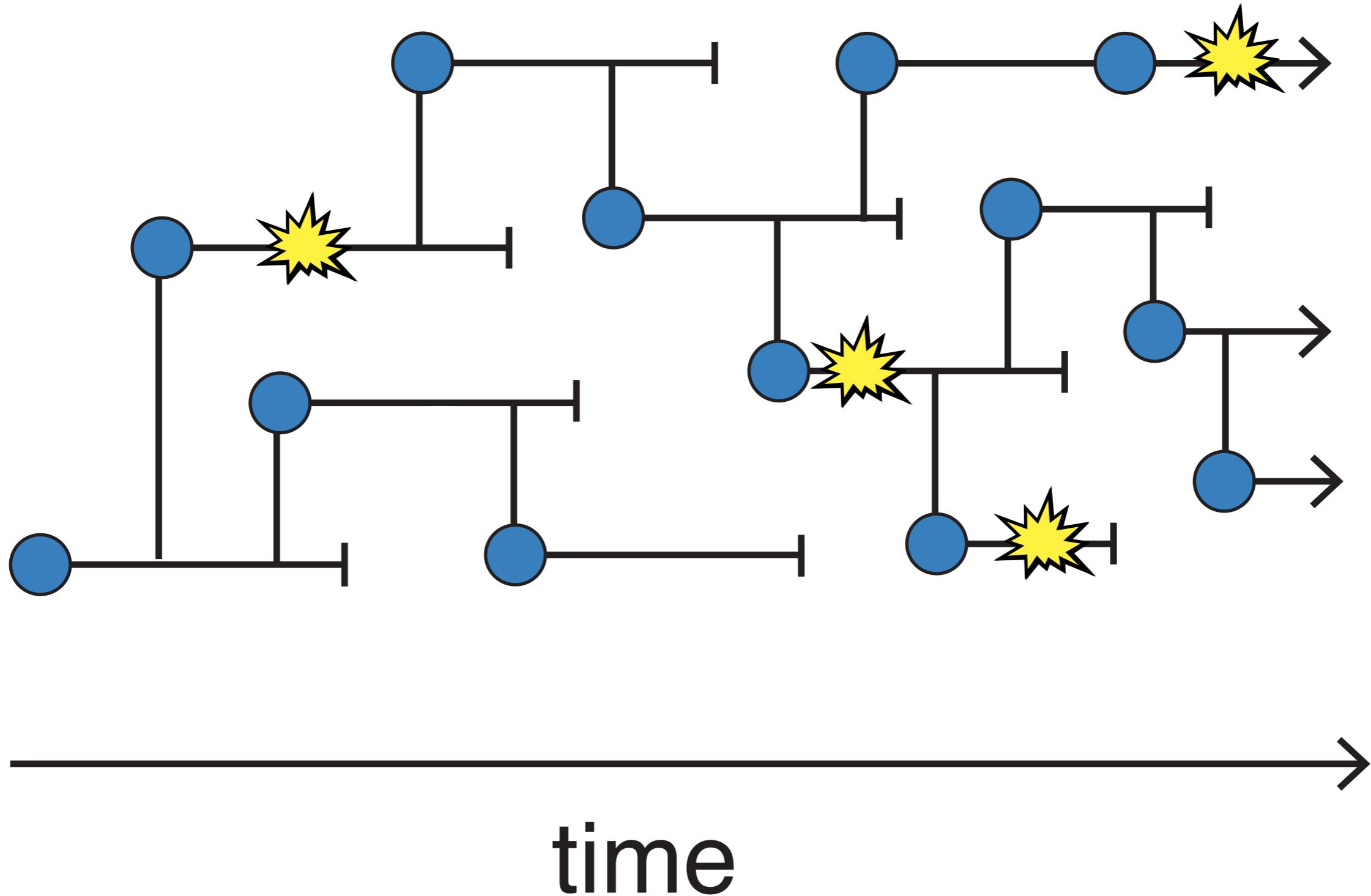
Viral genomes contain records of a virus's transmission history



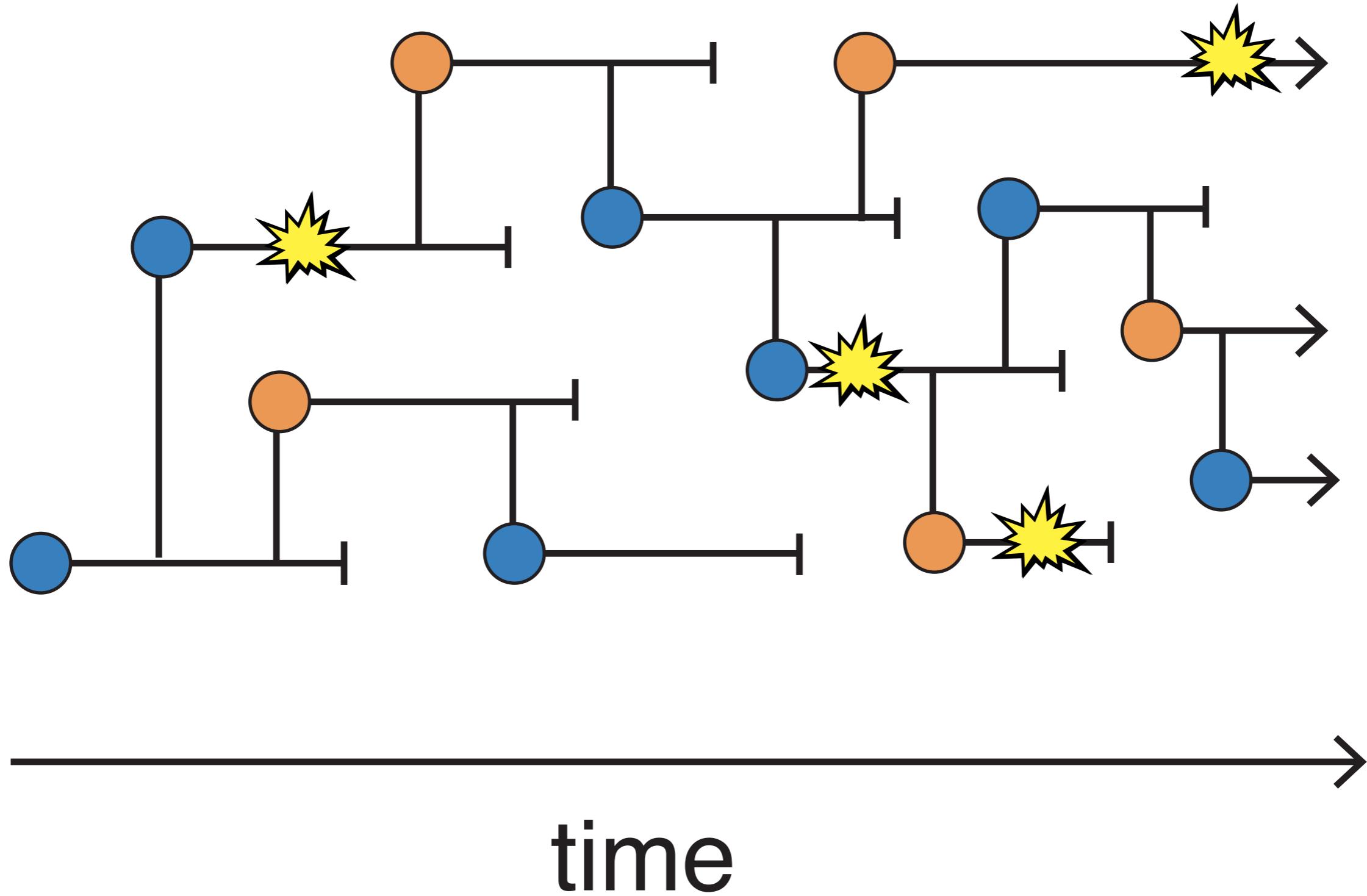
Viral genomes contain records of a virus's transmission history



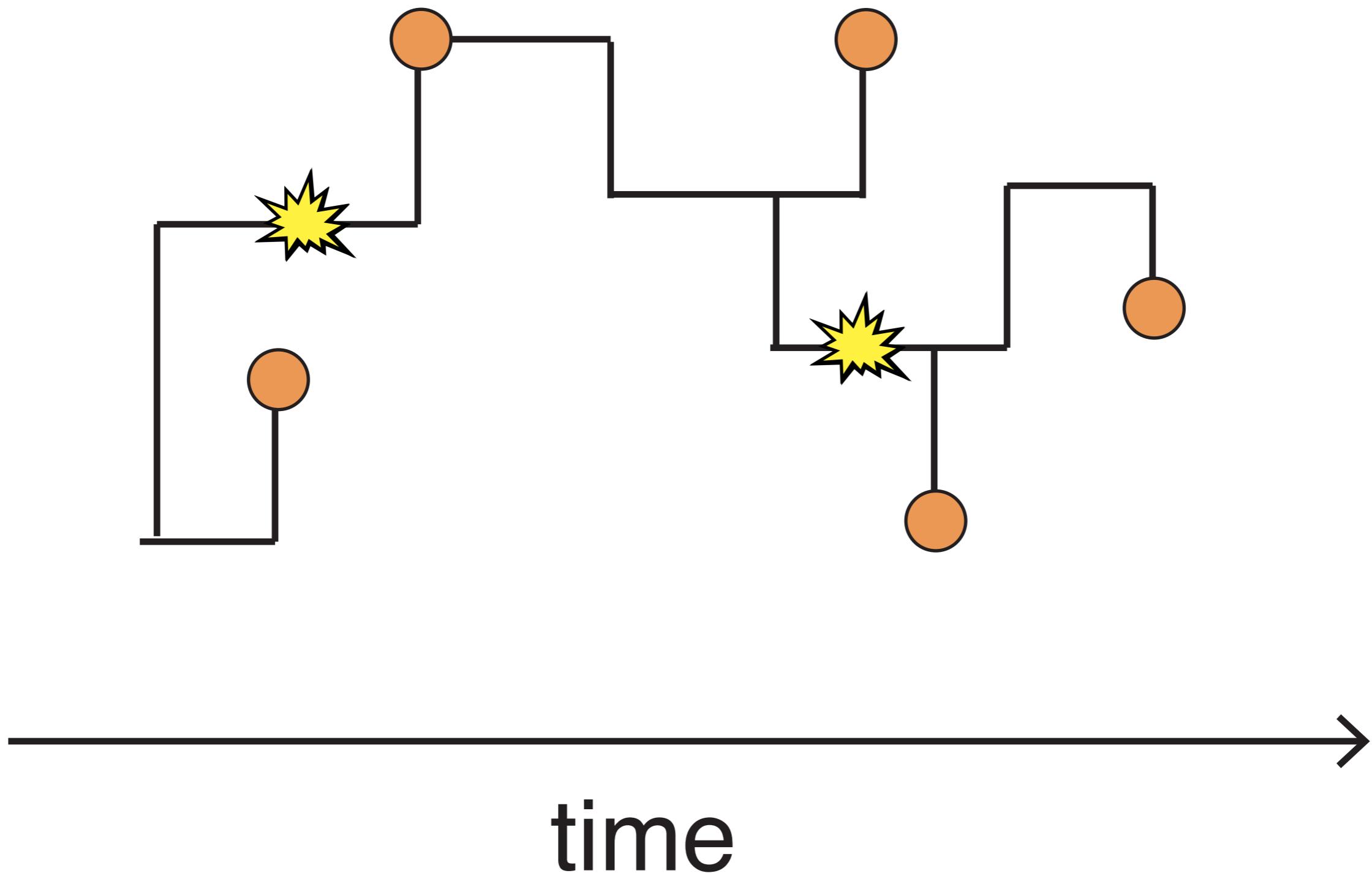
Viral genomes contain records of a virus's transmission history



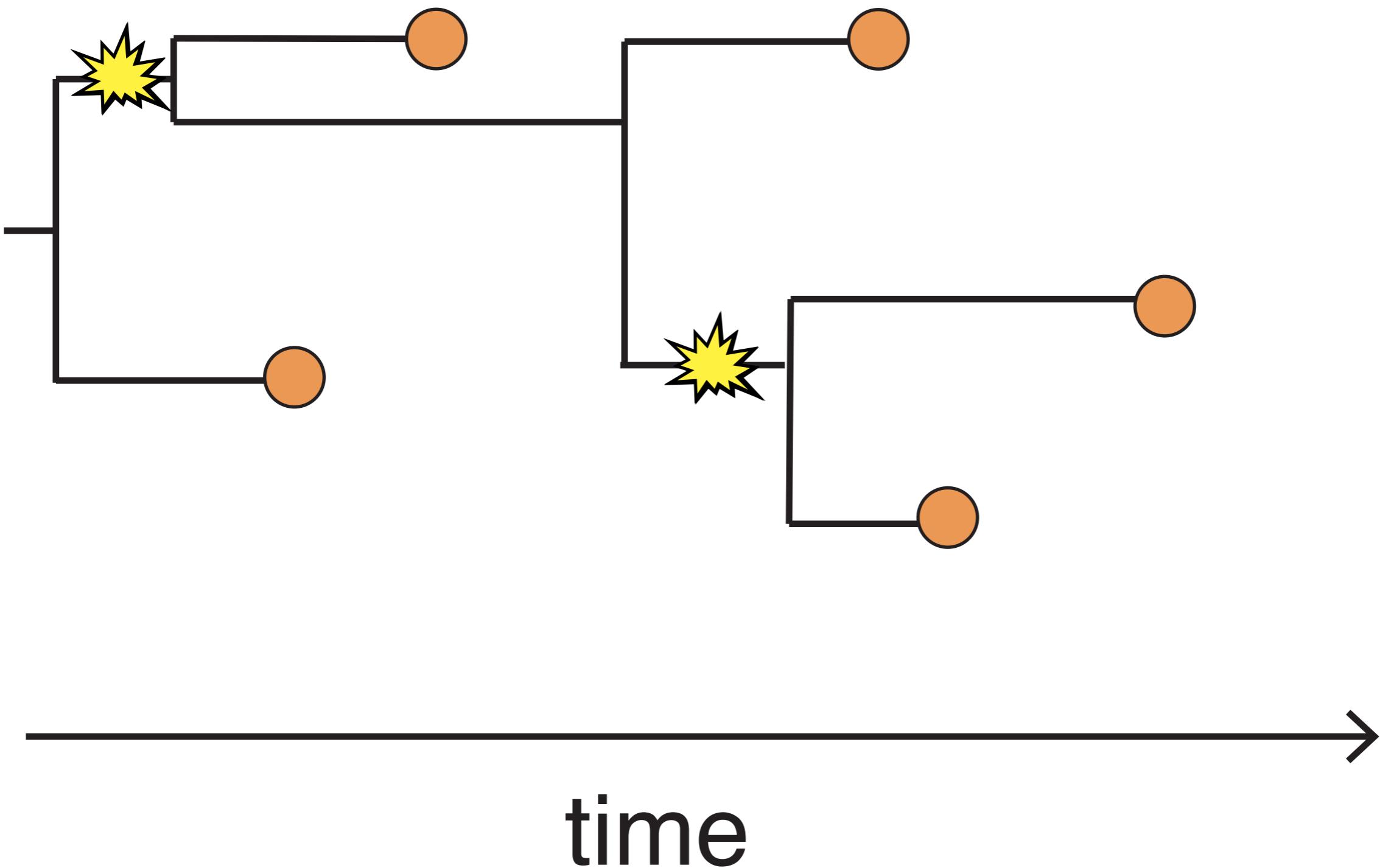
Viral genomes contain records of a virus's transmission history



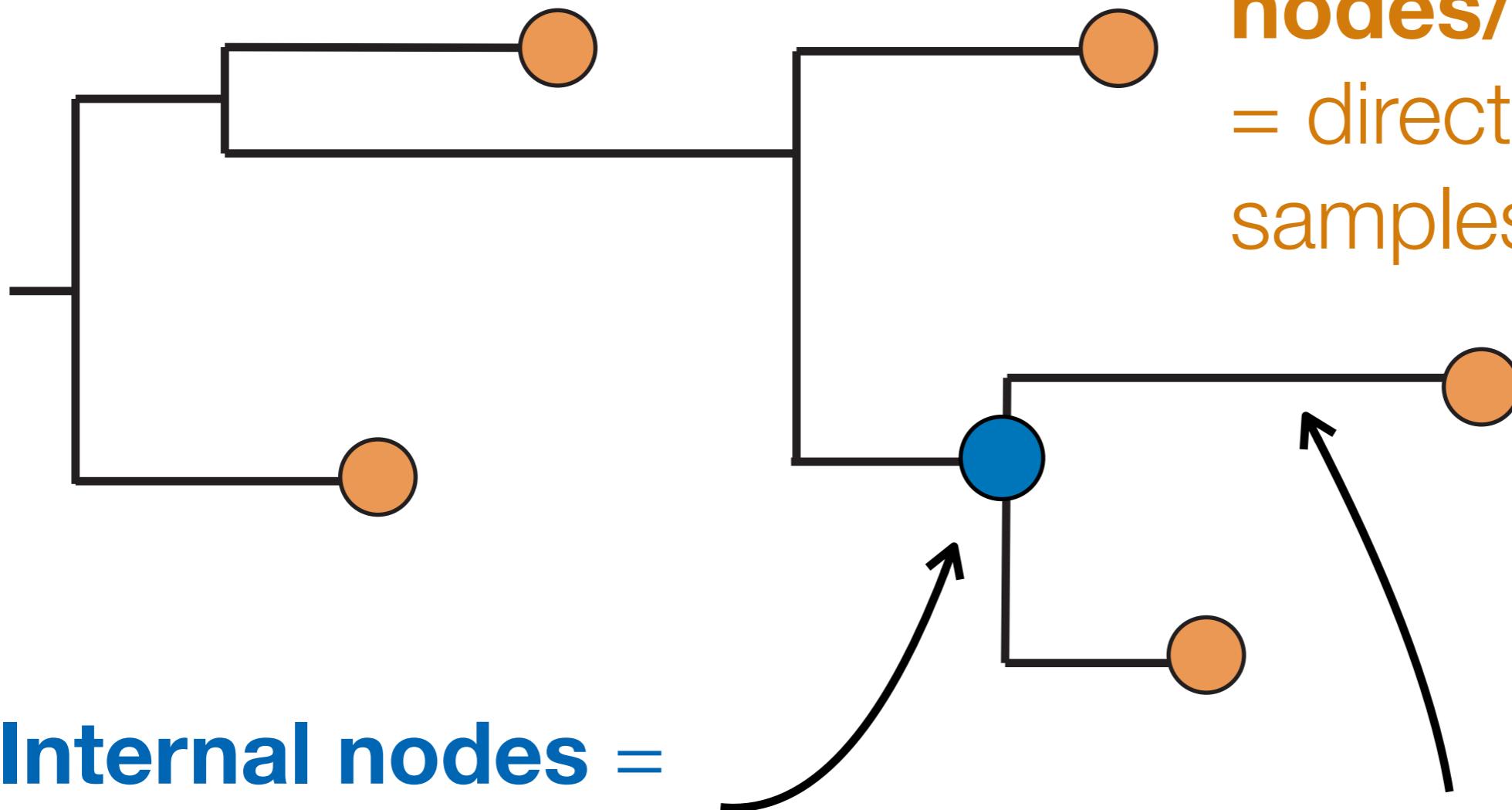
Viral genomes contain records of a virus's transmission history



Viral genomes contain records of a virus's transmission history



A phylogenetic tree is a hypothesis of shared descent

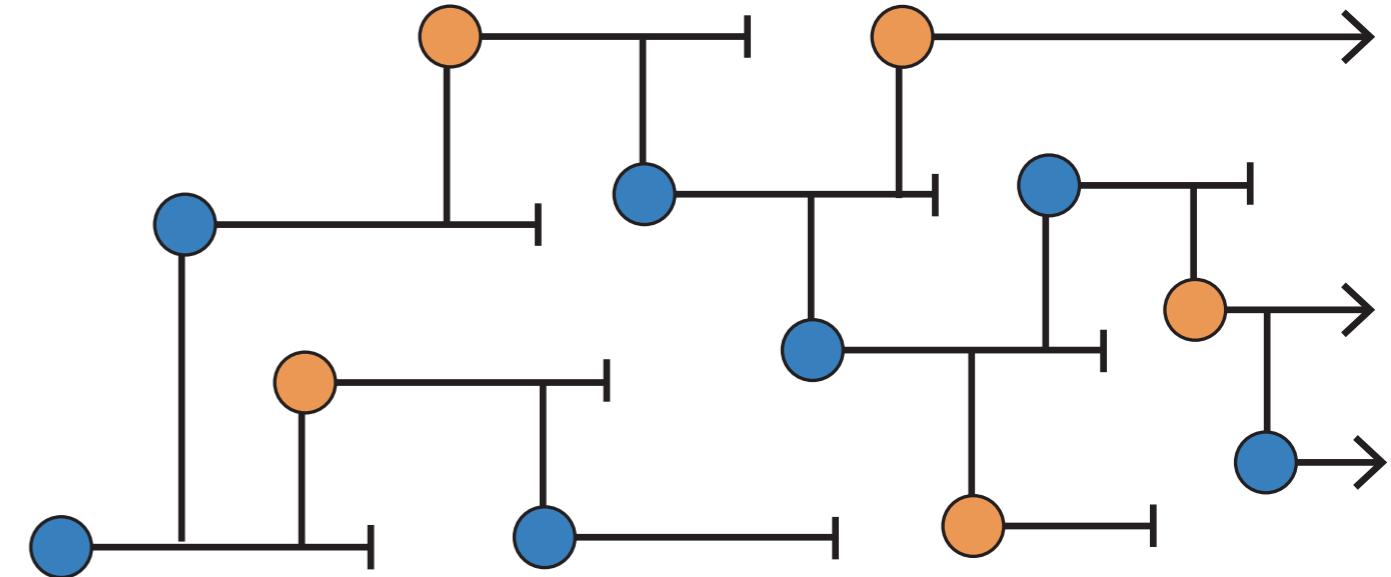
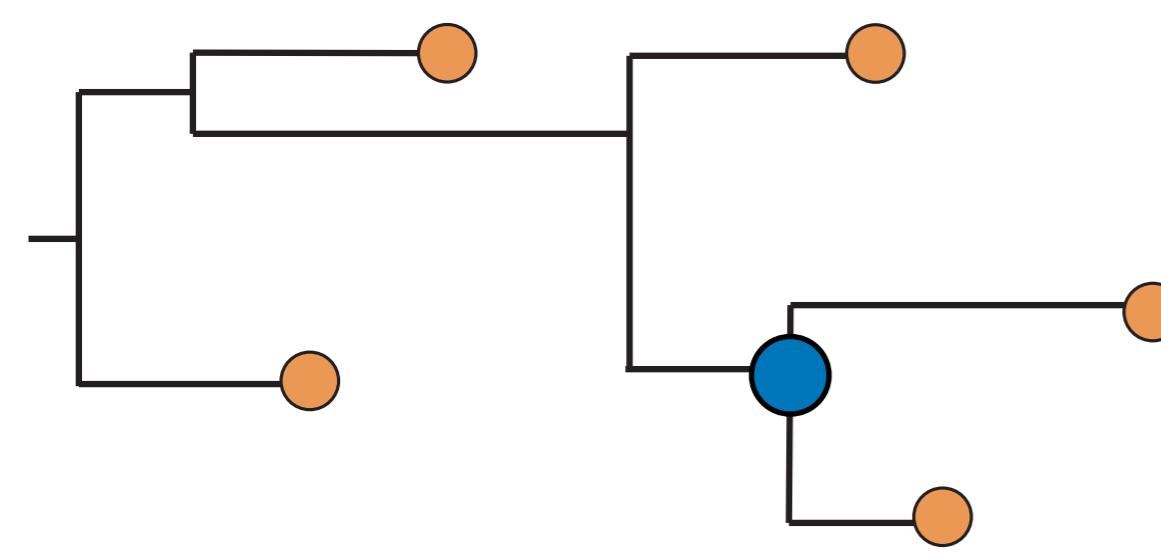


Internal nodes =
hypothetical,
inferred ancestors

**Tips/terminal
nodes/leaves/taxa**
= directly observed
samples

branches = a way to
connect tips and nodes; units
can be divergence or time

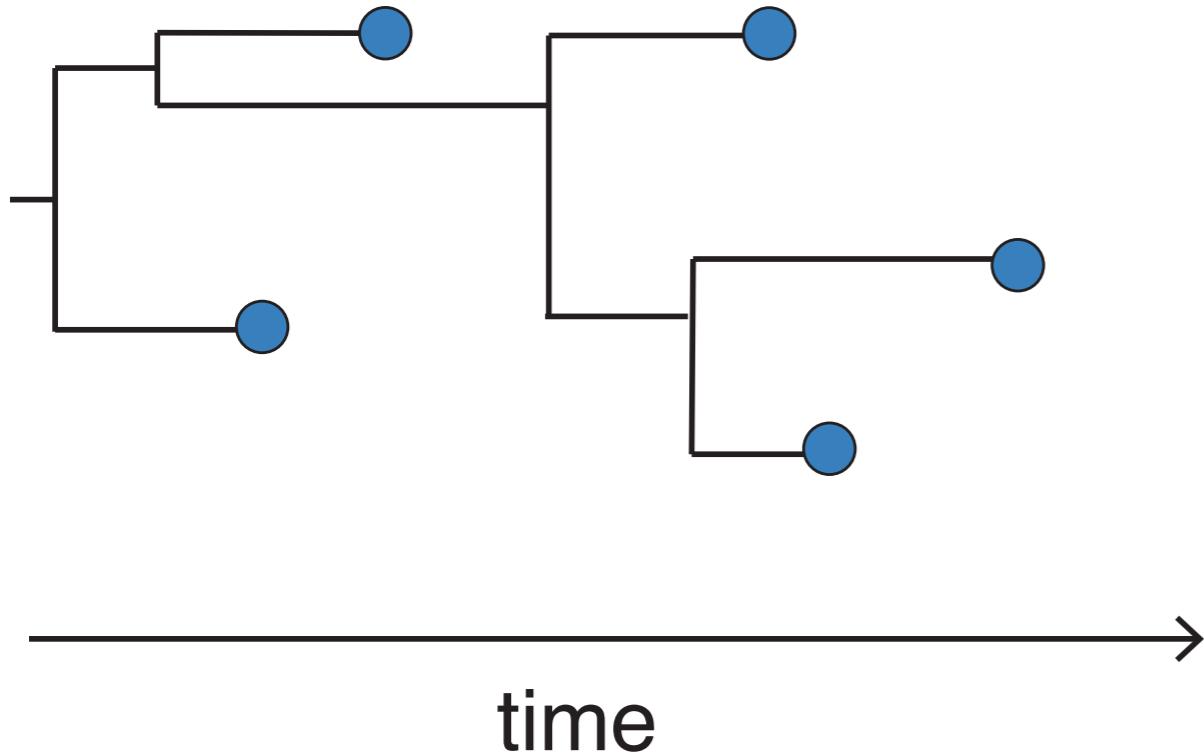
Phylogenetic trees are distinct from transmission trees



Phylogenetic tree:
hypothesis of shared
ancestry

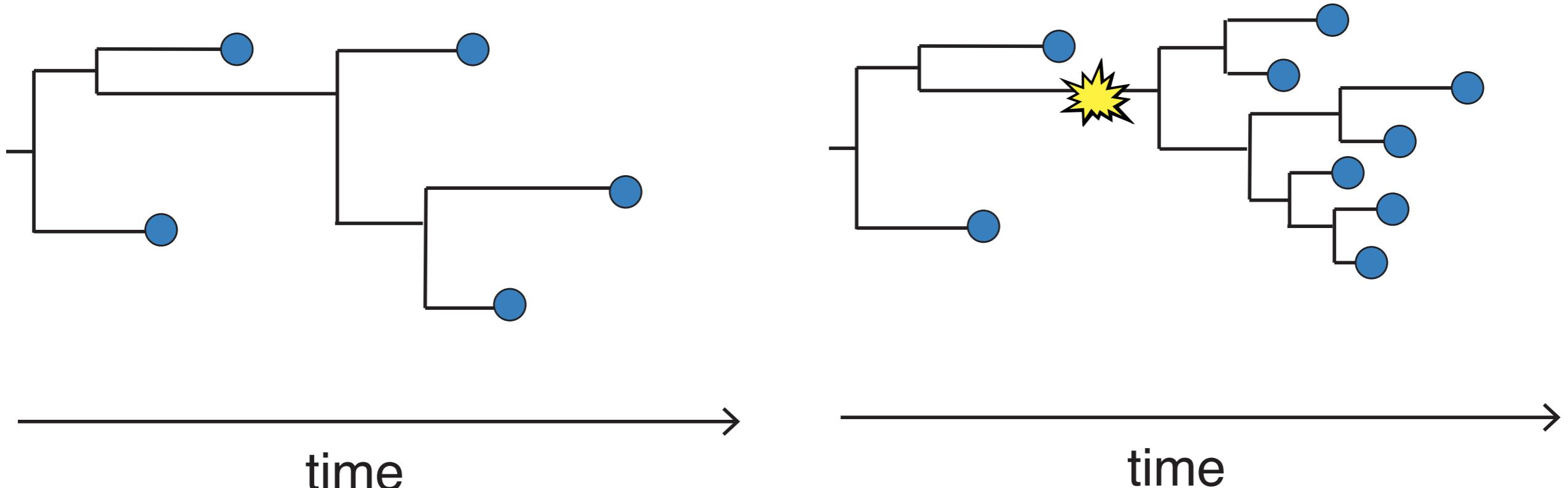
Transmission tree: true
transmission history

The shape of phylogenetic trees changes under different evolutionary, ecological, and epidemiologic scenarios



Original tree

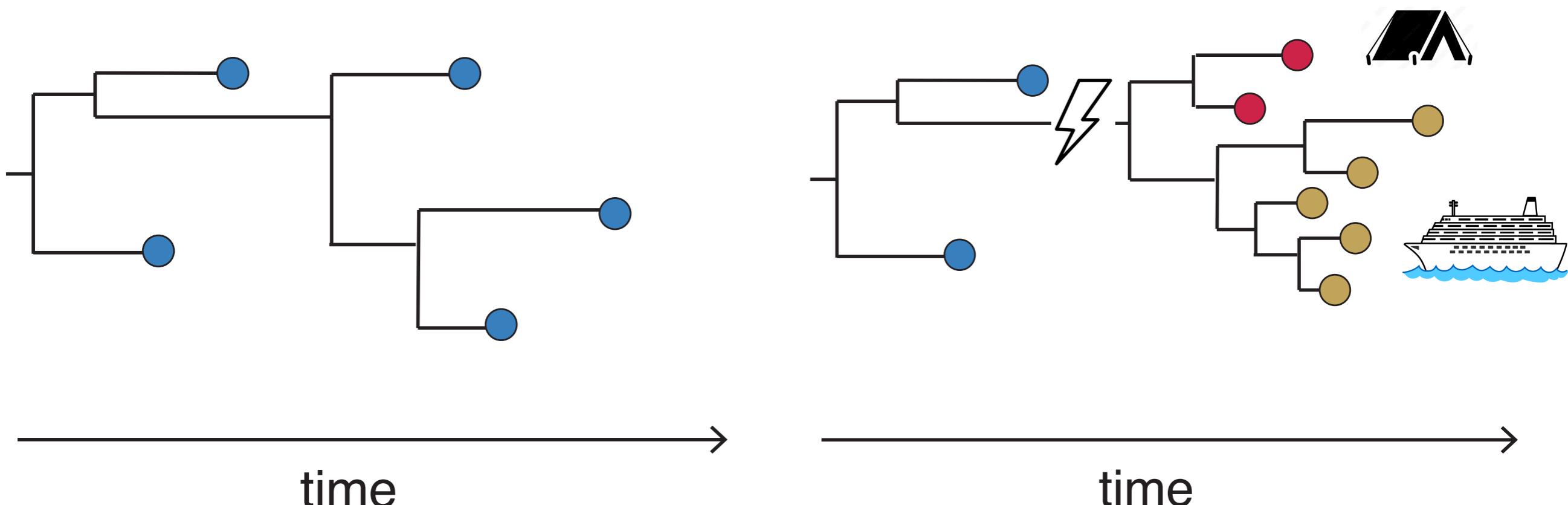
The shape of phylogenetic trees changes under different evolutionary, ecological, and epidemiologic scenarios



Original tree

Increased
transmissibility

The shape of phylogenetic trees changes under different evolutionary, ecological, and epidemiologic scenarios



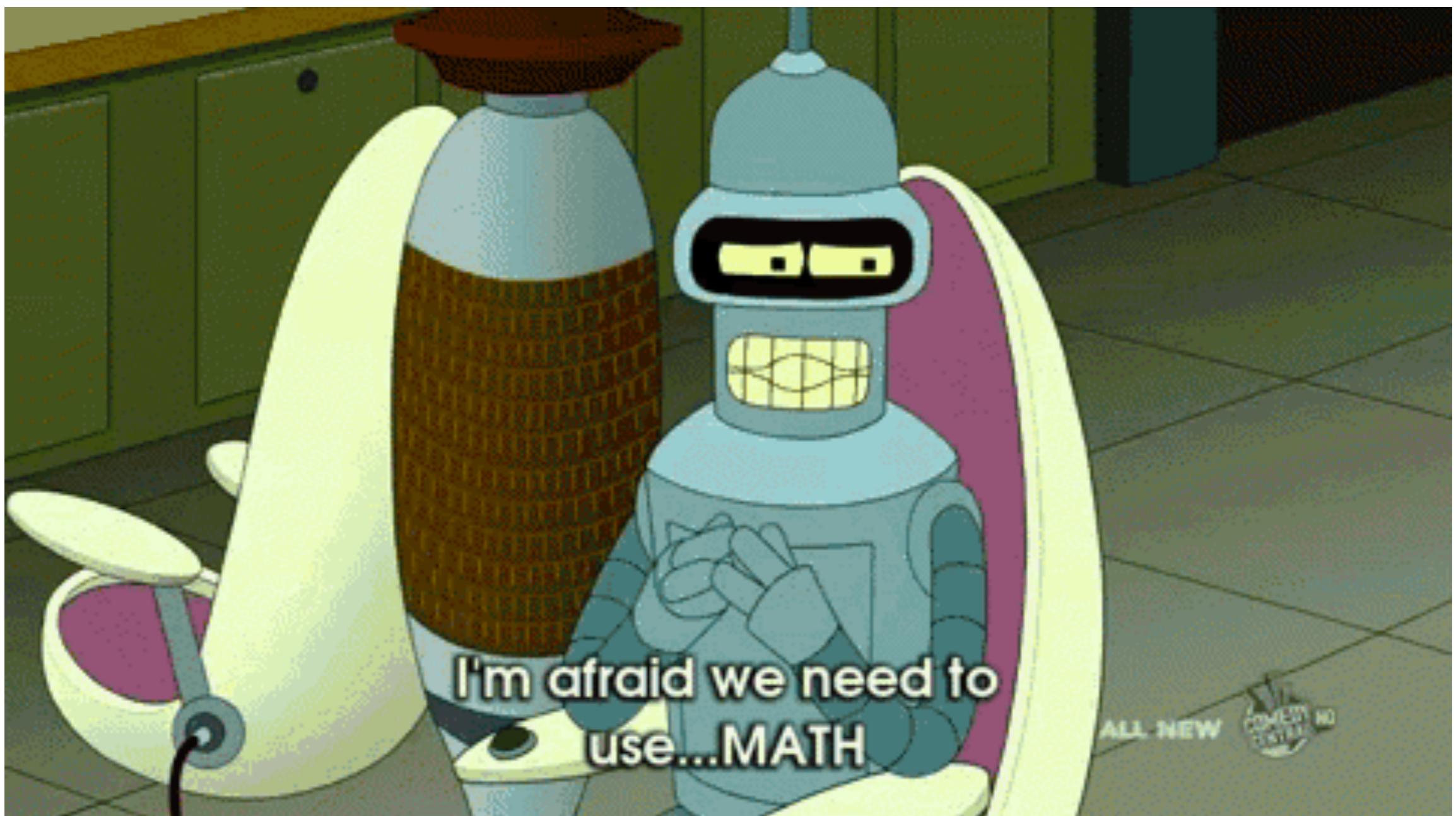
Original tree

Population
subdivision

... but there's a problem

We don't usually know the evolutionary or transmission history of a pathogen.

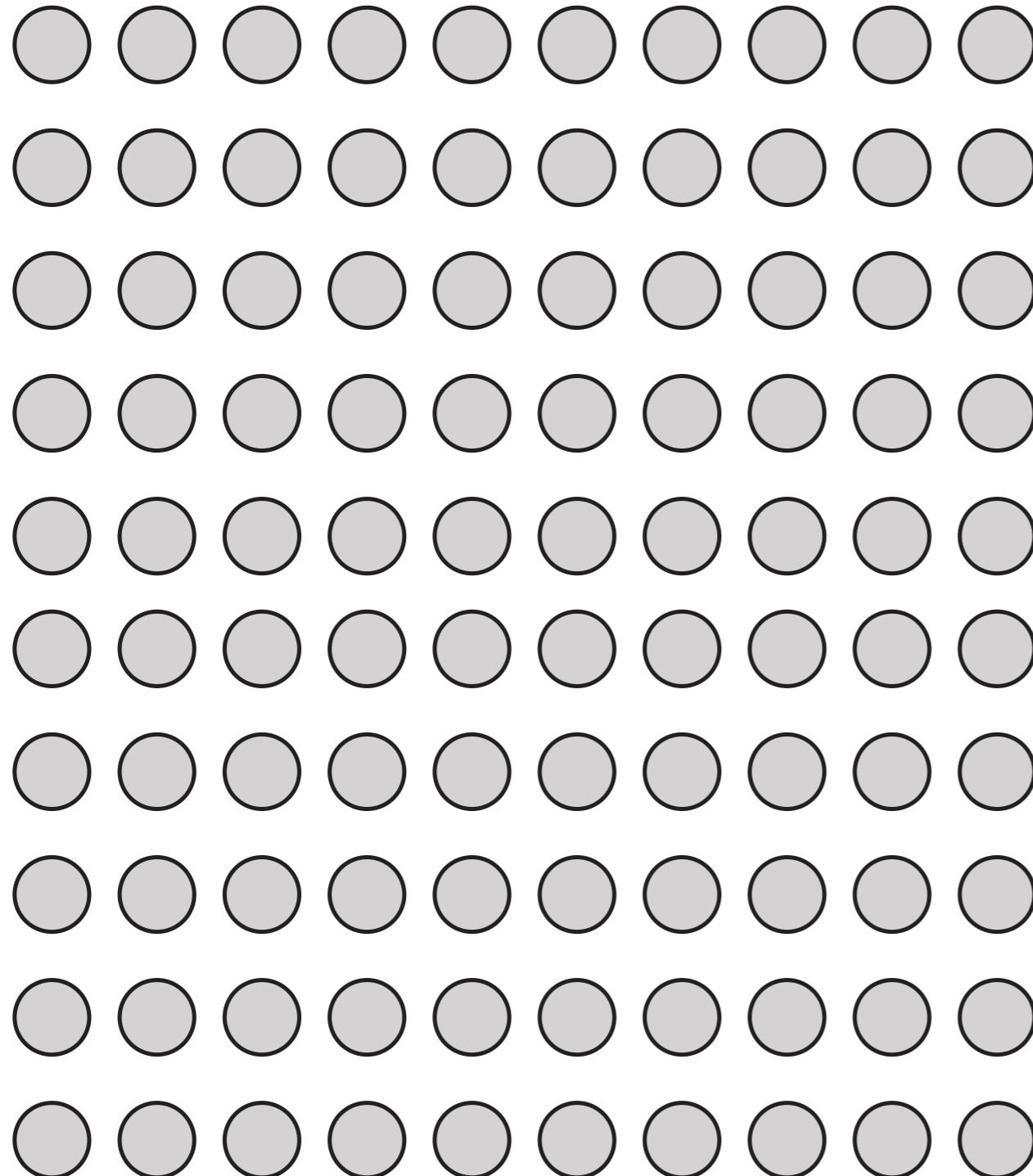
We need a way to quantitatively infer these histories, given a phylogeny



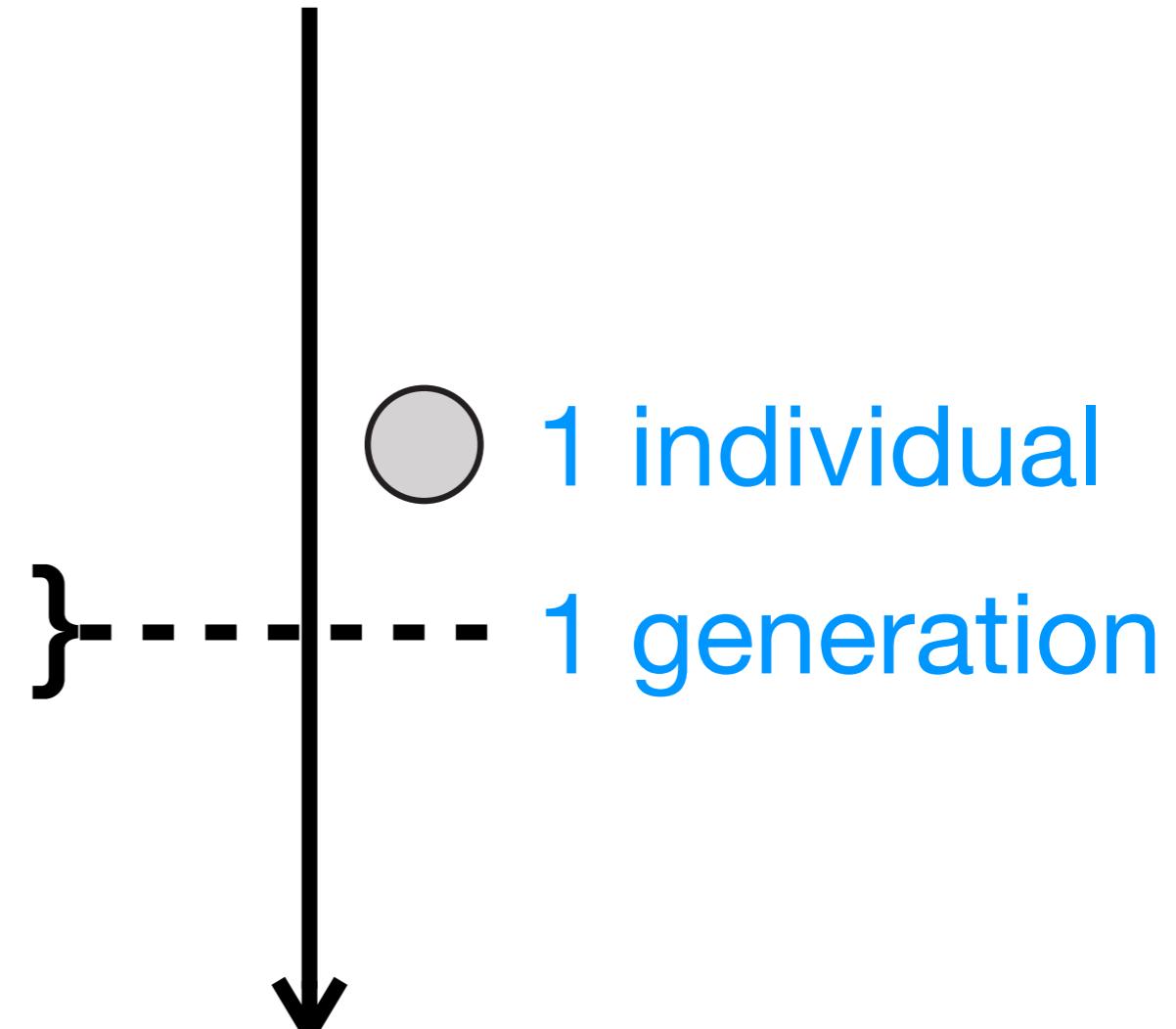
The coalescent is a mathematical model that describes the ancestry of a sample of non recombining gene copies.

* Volz, Koelle, and Bedford, PLOS Computational Biology, 2013

The coalescent

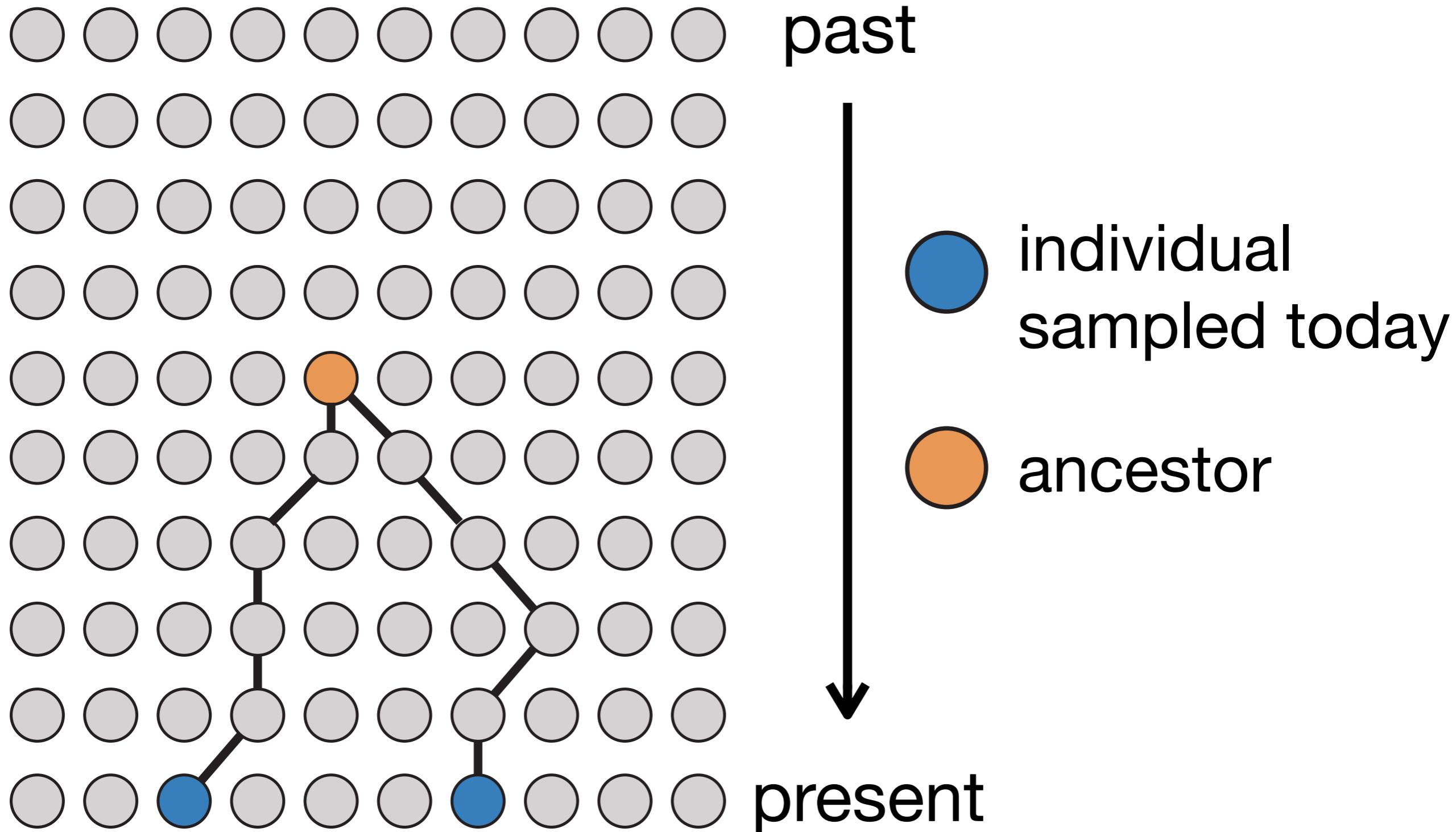


the past

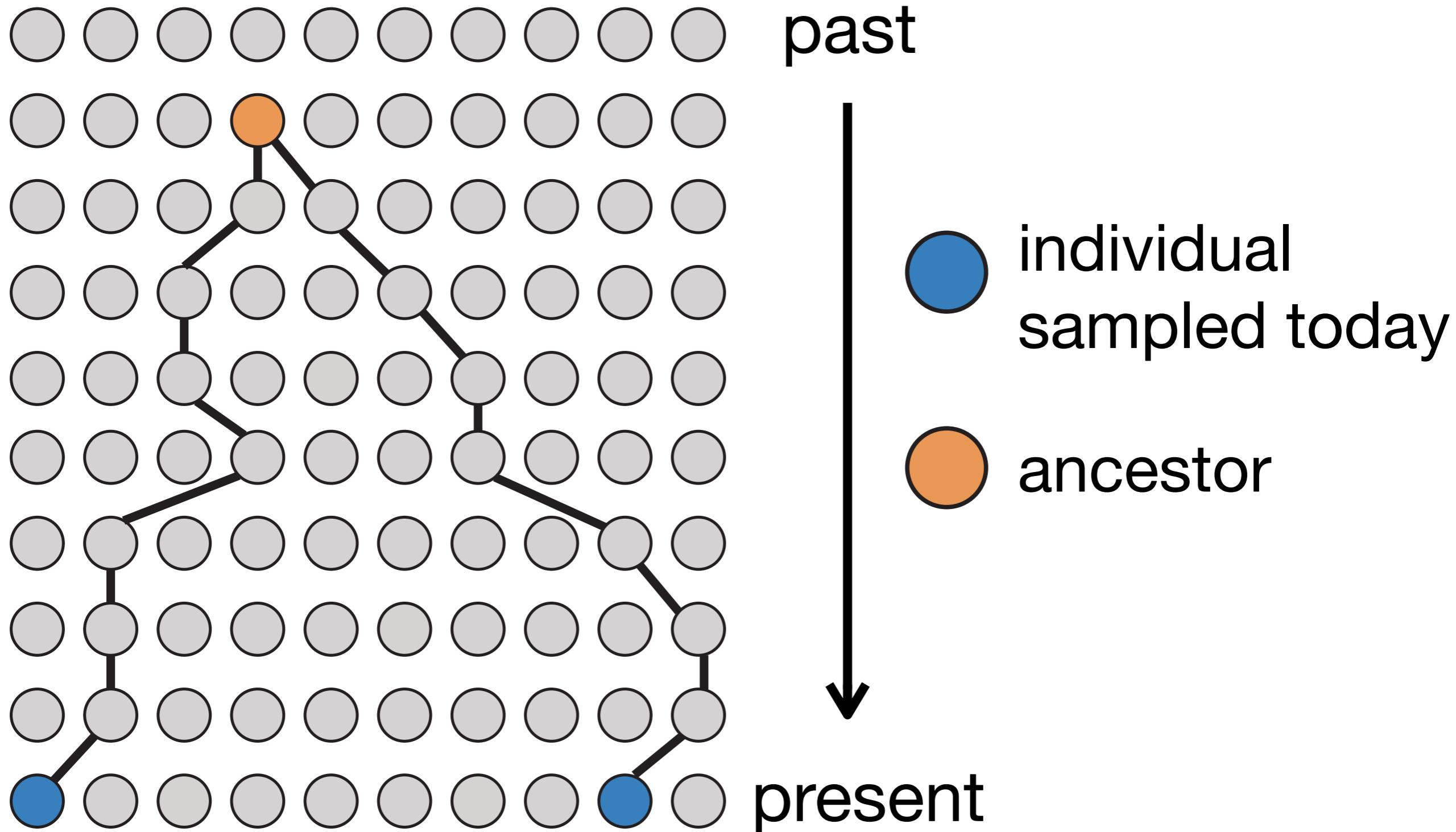


present day

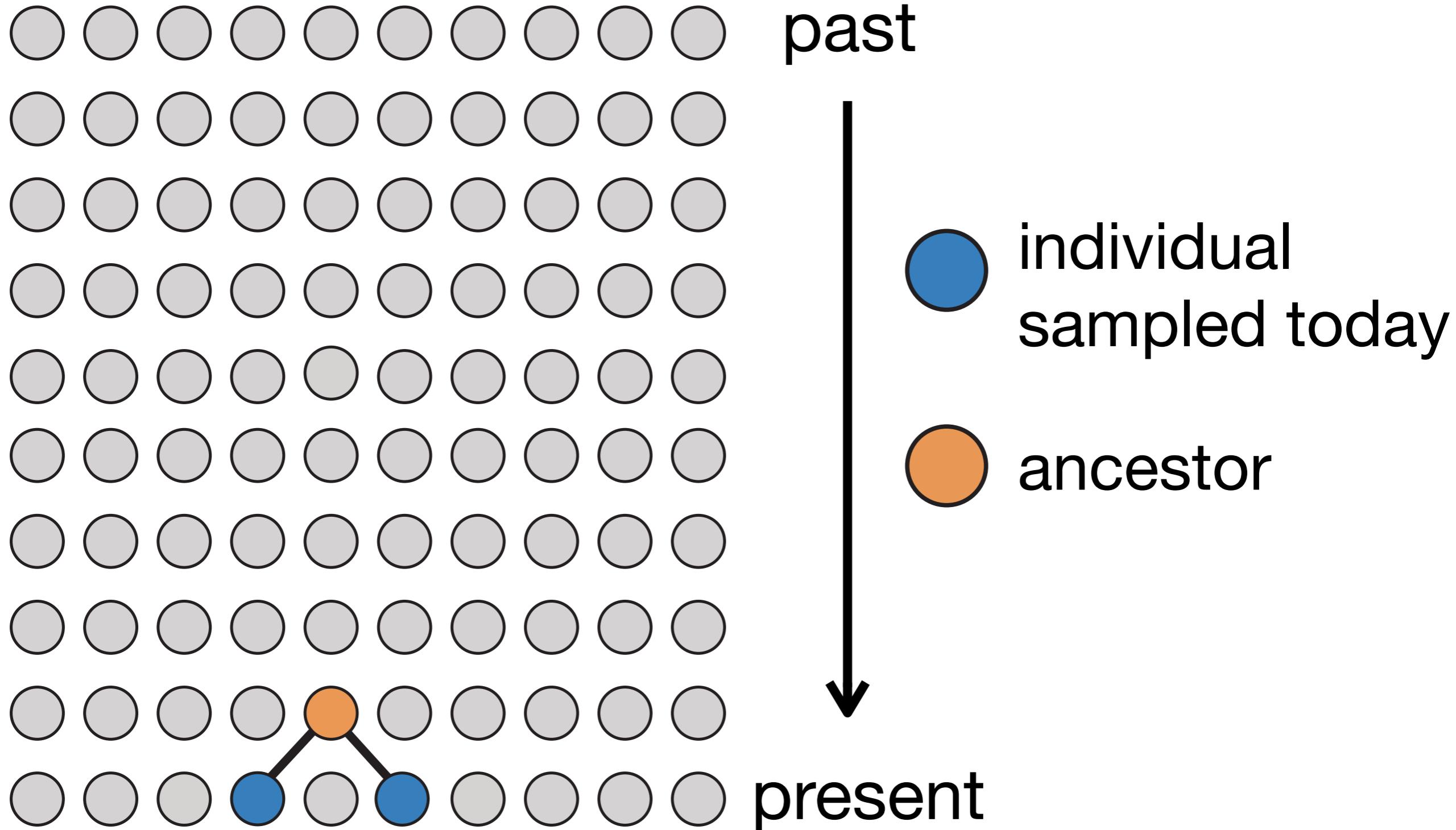
Any 2 individuals sampled today will share a common ancestor in the past



Any 2 individuals sampled today will share a common ancestor in the past

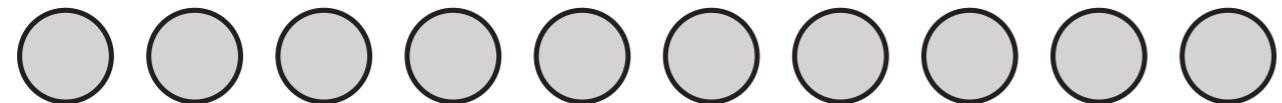


Any 2 individuals sampled today will share a common ancestor in the past



The expected time to coalescence depends on population size

past



past



past



past



past



past



past



past



past



past



present



In any given generation, the probability that 2 individuals coalesce = $1/N$
(N = population size)

The expected time to coalescence depends on population size

In any given generation, the **probability** that 2 individuals pick the same parent and coalesce = $1/N$

The expected **waiting time** for 2 lineages to coalesce is geometrically distributed with mean = N

* N is in units of generations

How can we scale this to more than 2 lineages?

In any given generation, the **probability** that i individuals pick the same parent and coalesce:

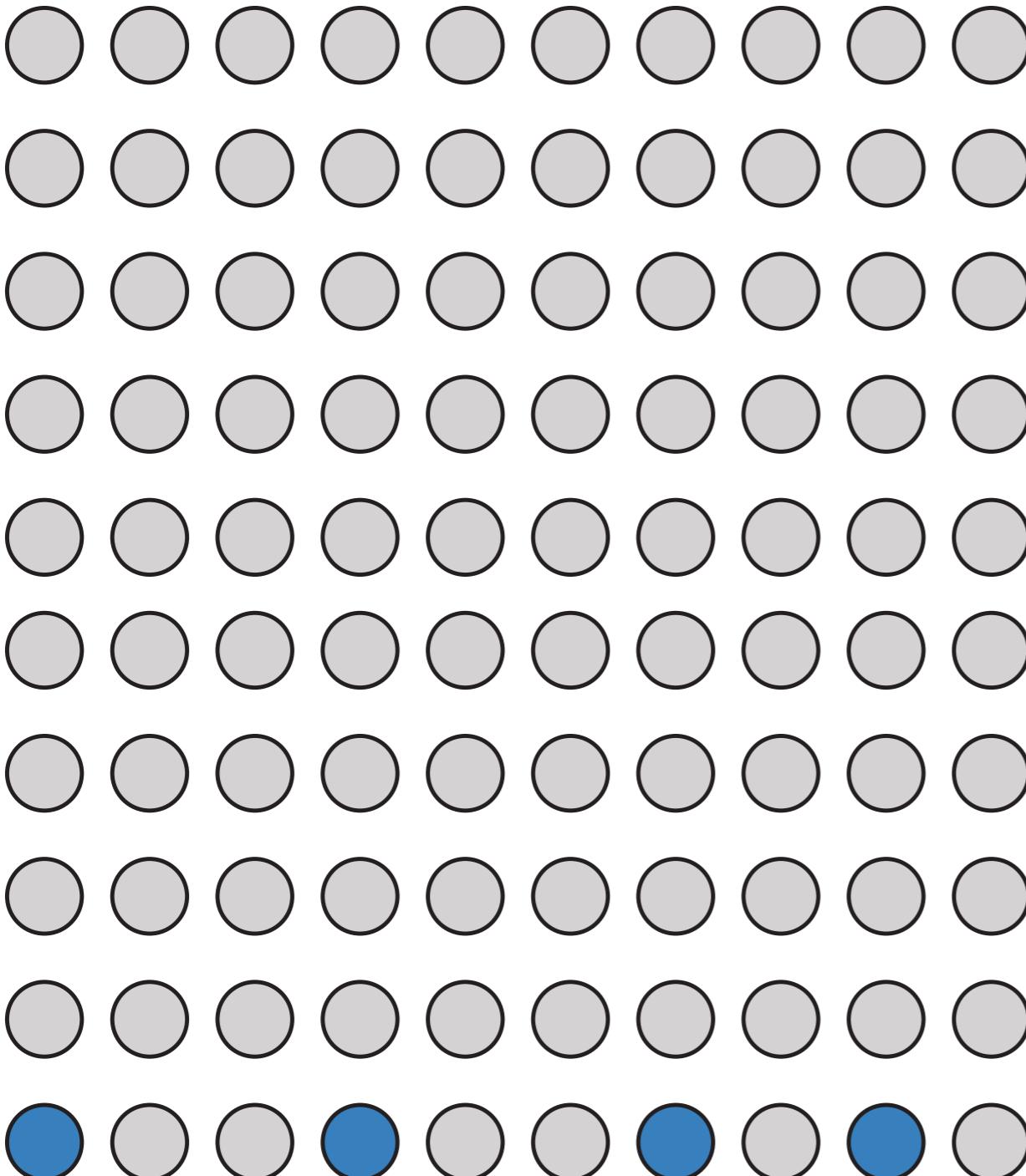
$$\frac{i(i-1)}{2N}$$

The expected **waiting time** i lineages to coalesce is exponentially distributed with a mean of:

$$\frac{2N}{i(i-1)}$$

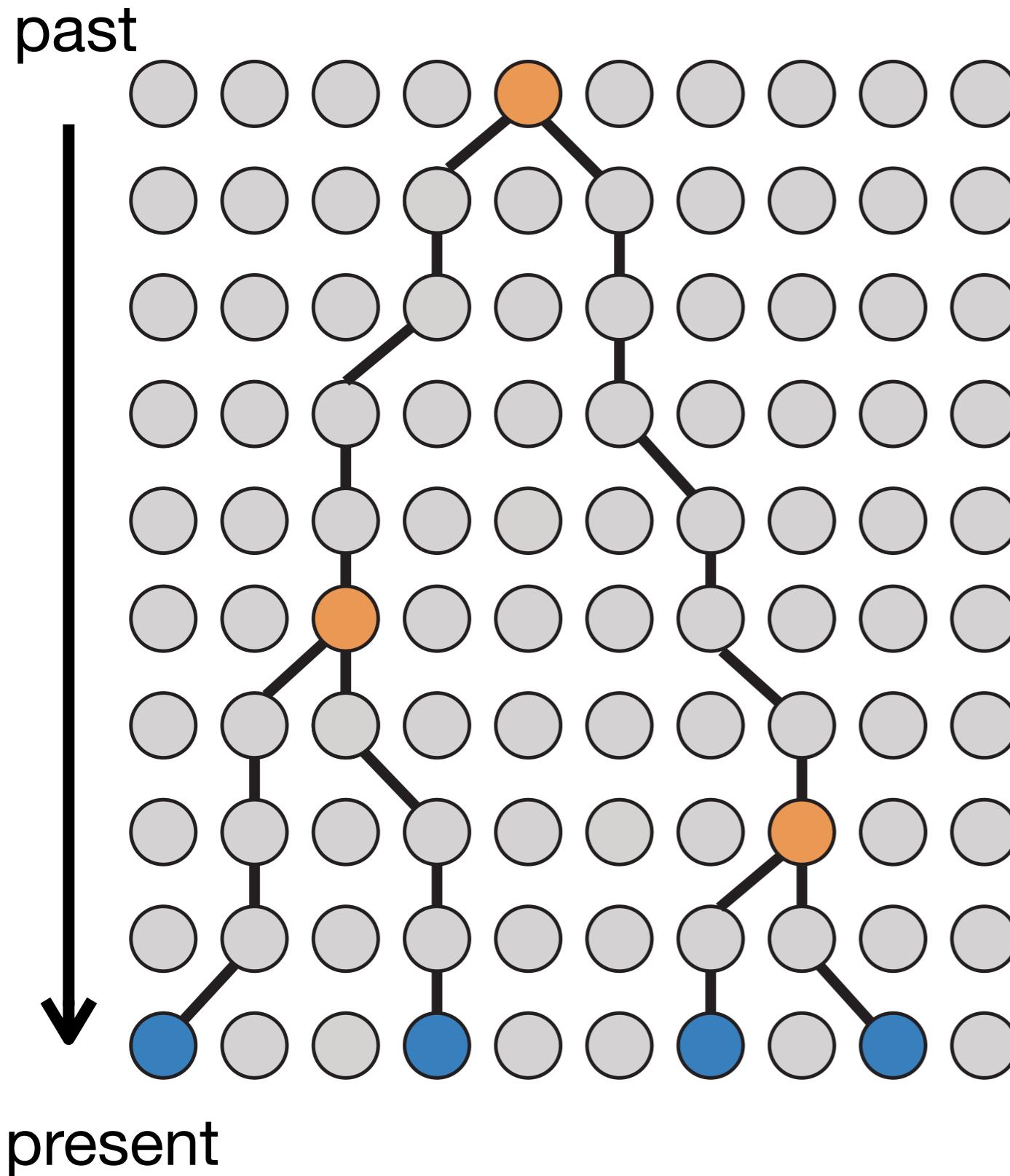
Coalescent times depend on population size and number of “lineages”

past

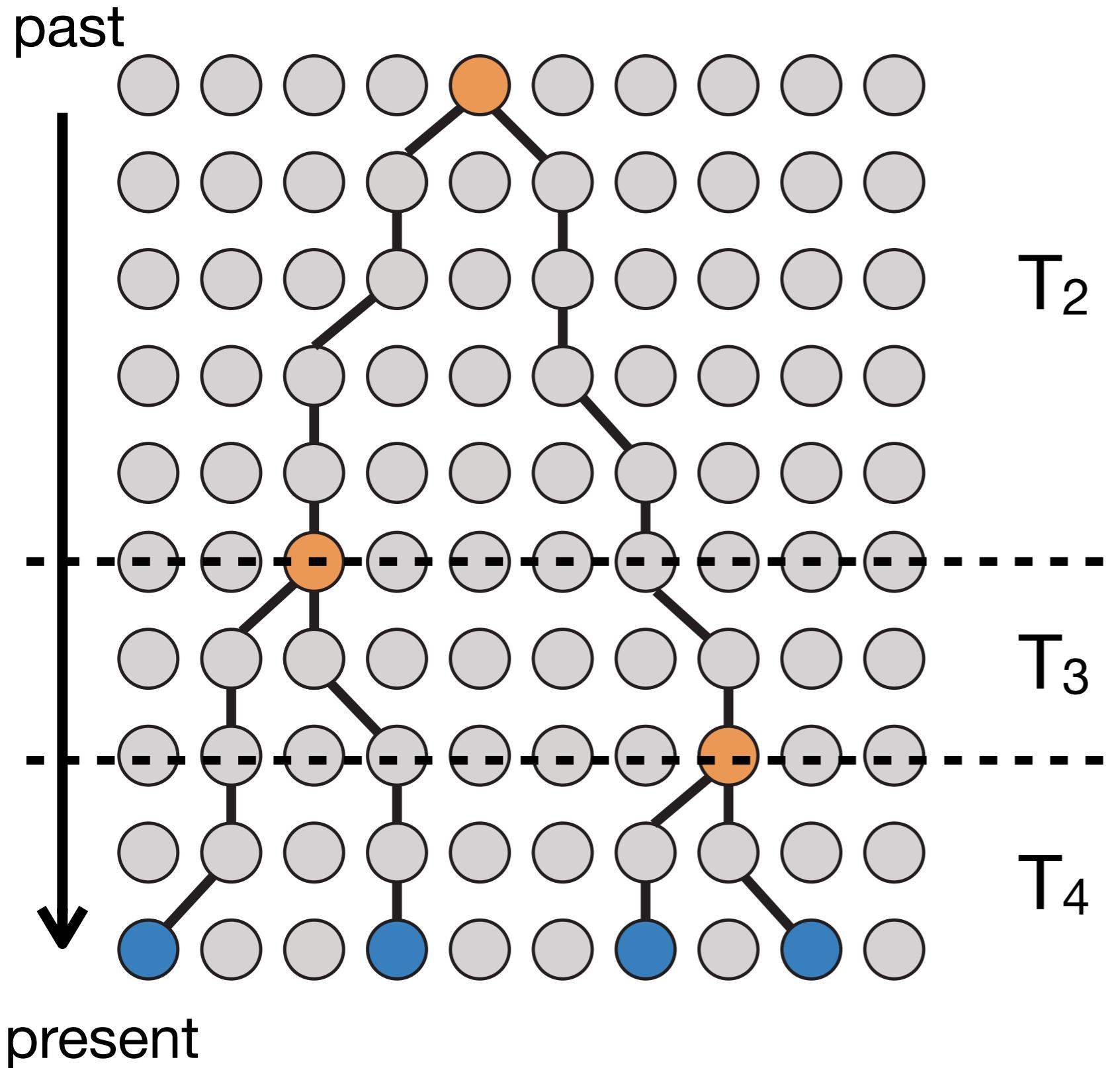


present

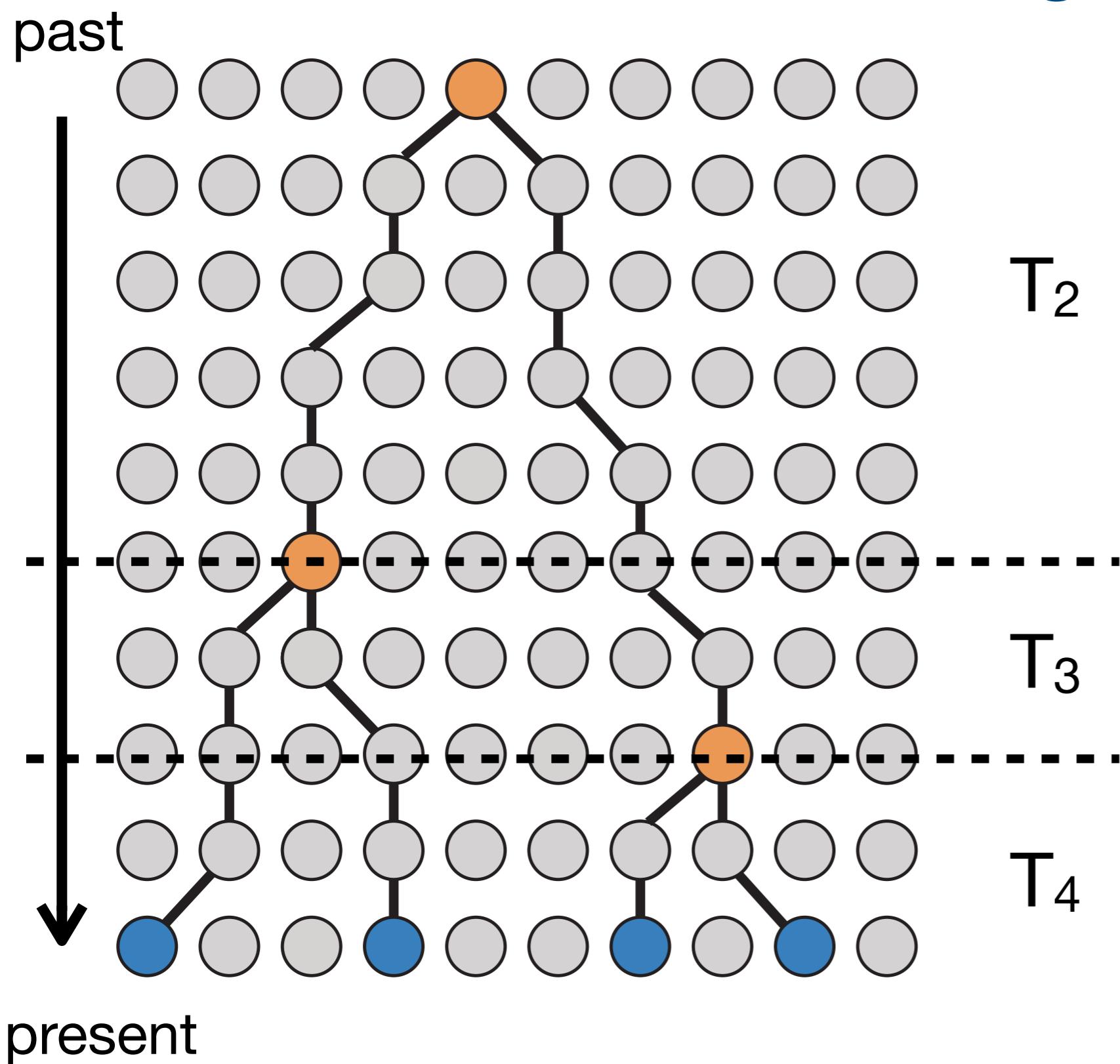
Coalescent times depend on population size and number of “lineages”



Coalescent times depend on population size and number of “lineages”



Coalescent times depend on population size and number of “lineages”

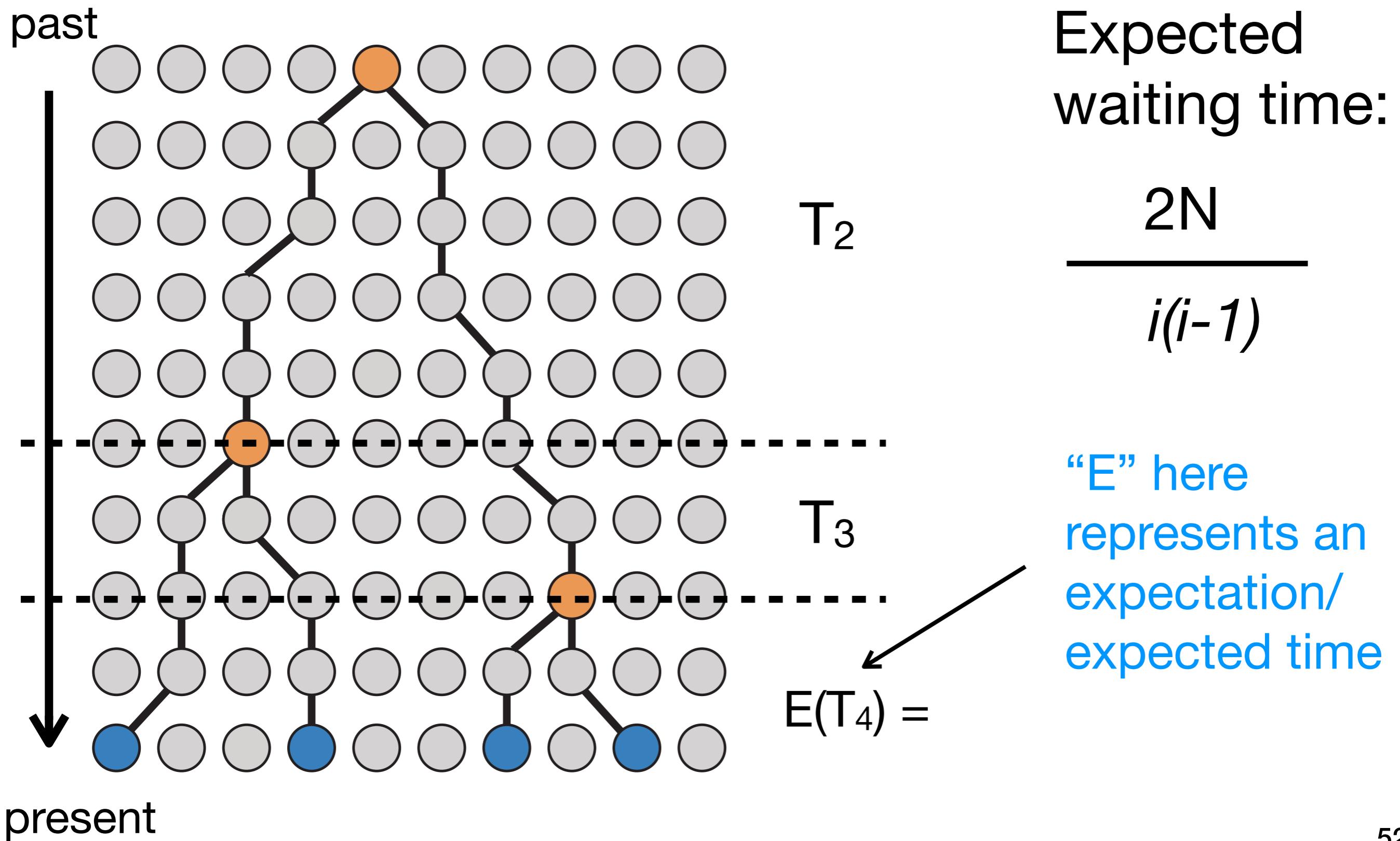


Expected waiting time:

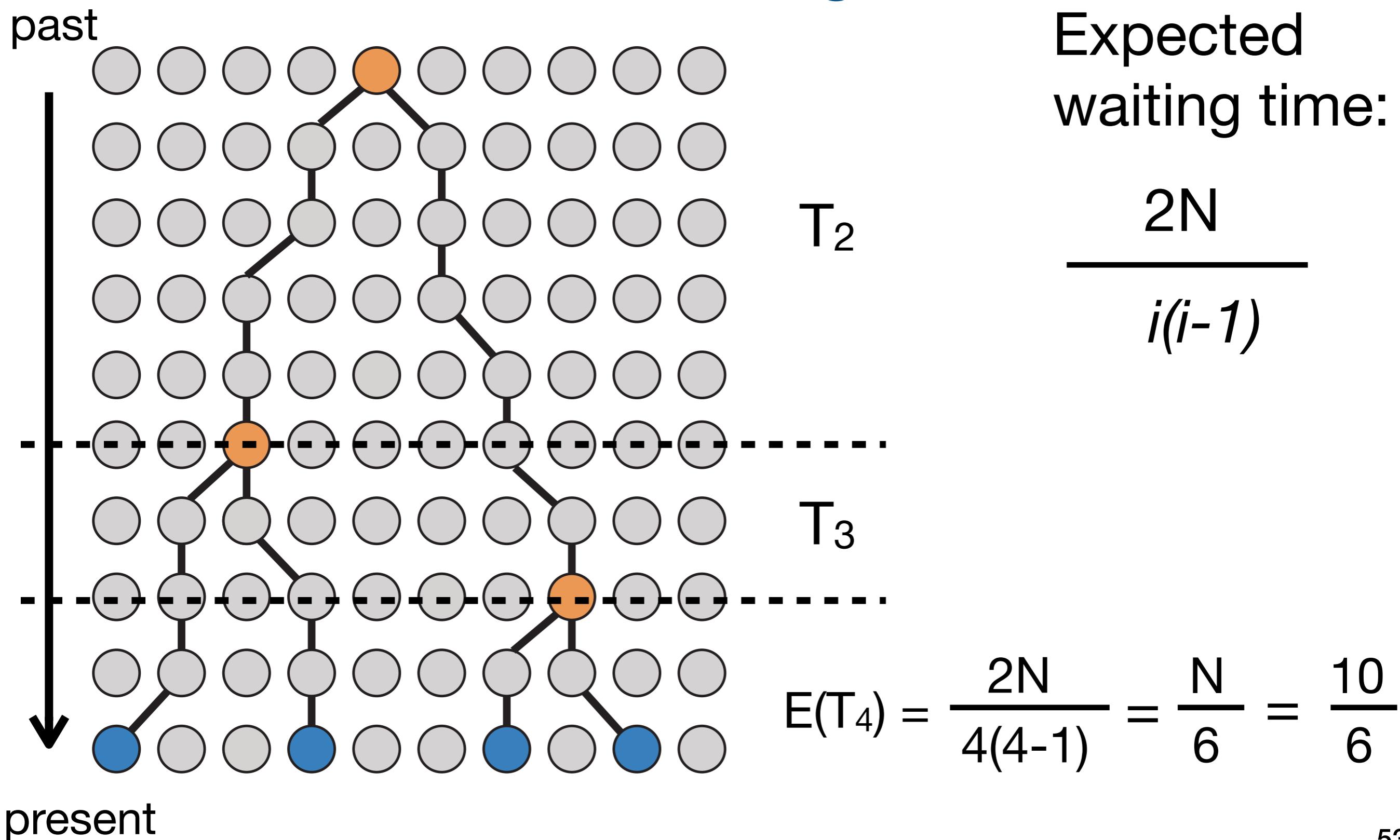
$$\frac{2N}{i(i-1)}$$

i = # of lineages
 N = population size, scaled by # of generations

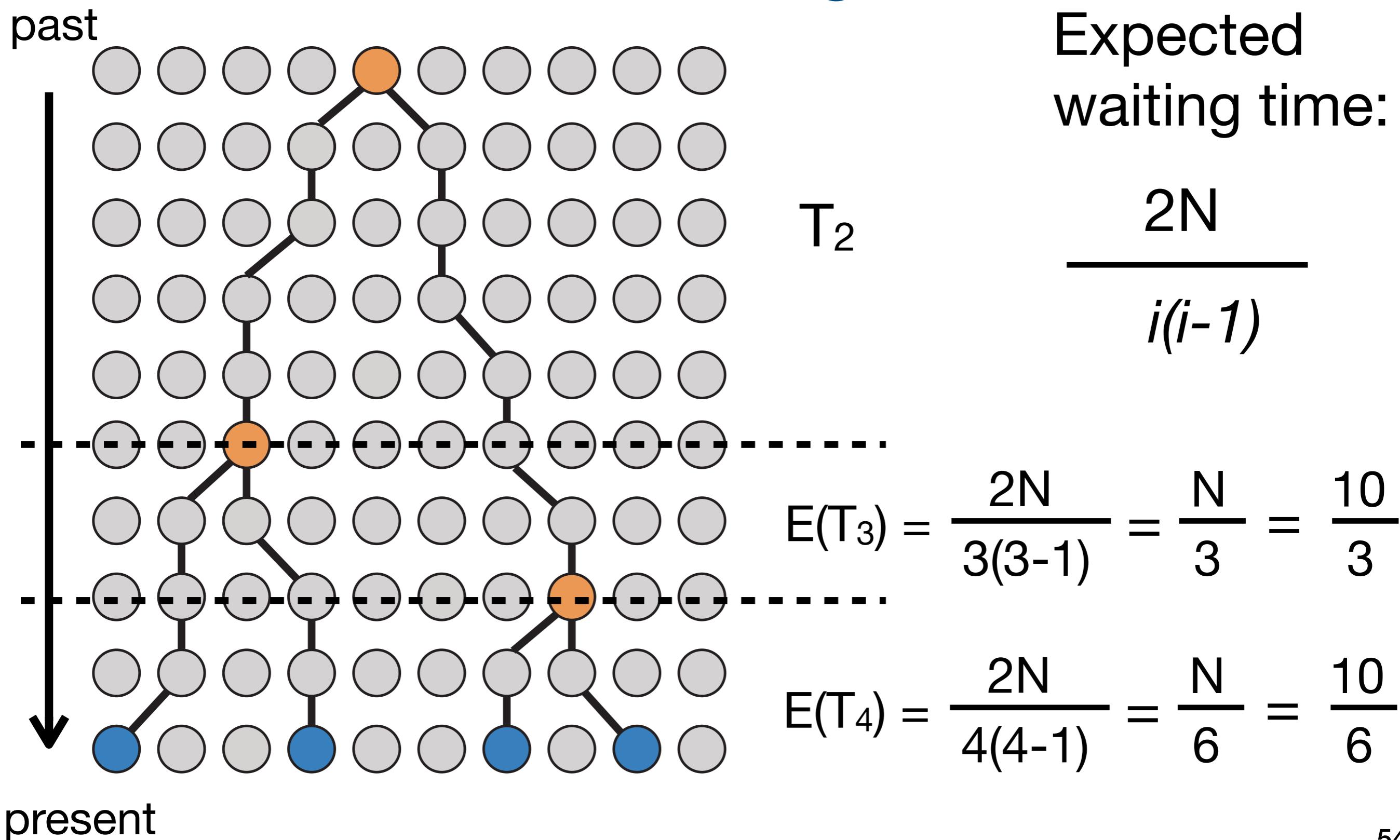
Coalescent times depend on population size and number of “lineages”



Coalescent times depend on population size and number of “lineages”

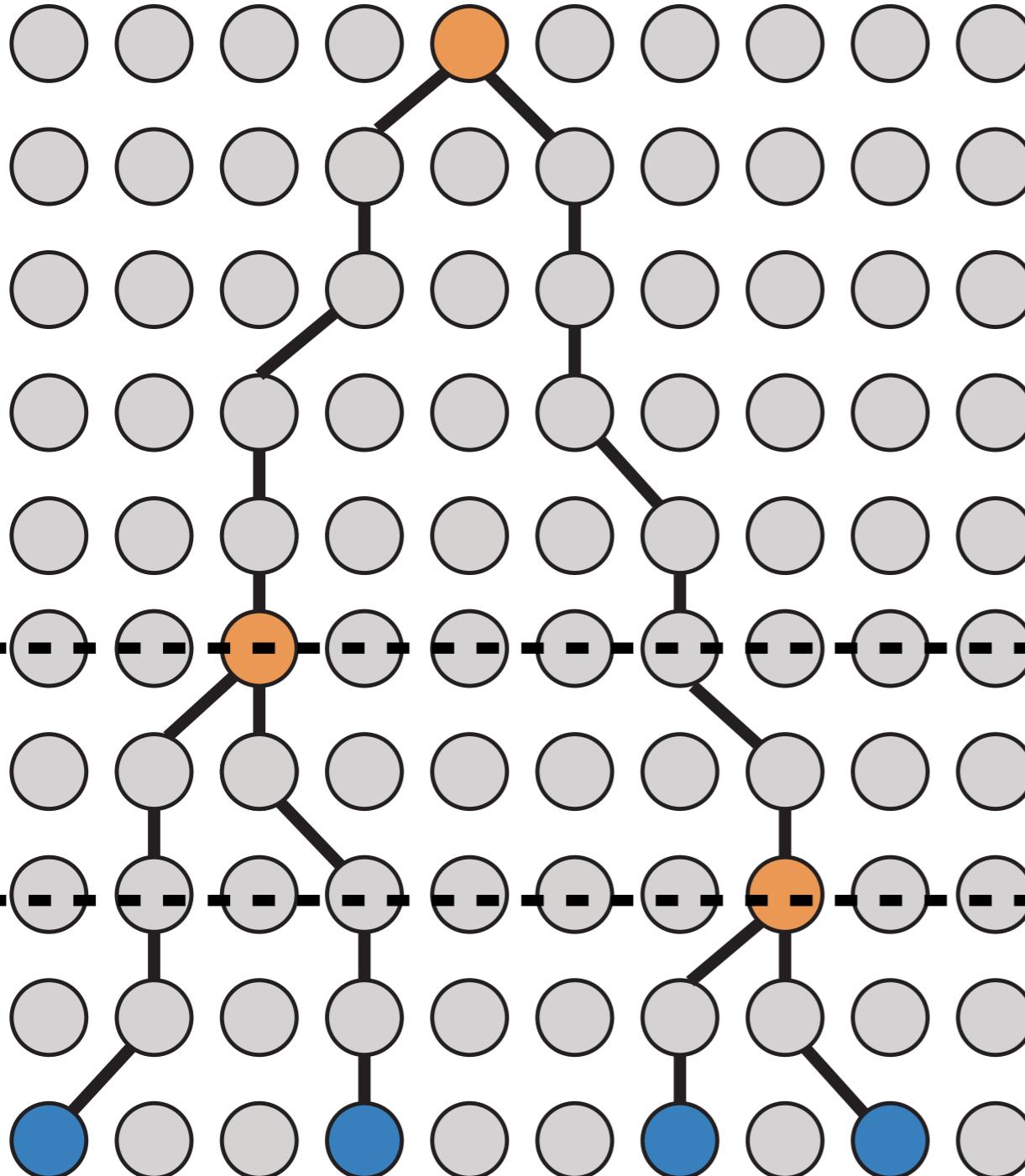


Coalescent times depend on population size and number of “lineages”



Coalescent times depend on population size and number of “lineages”

past



$$E(T_2) = \frac{2N}{2(2-1)} = N = 10$$

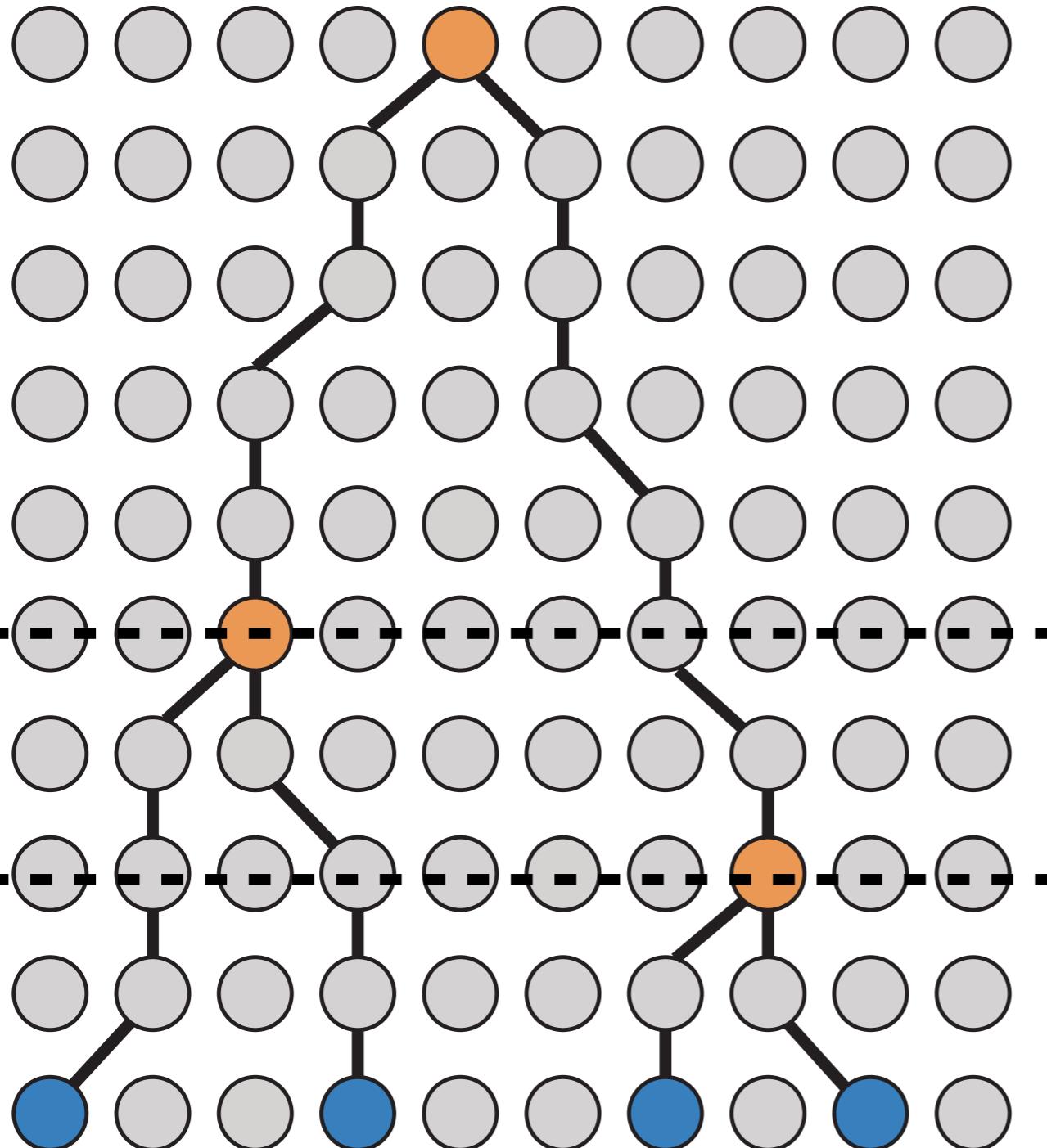
$$E(T_3) = \frac{2N}{3(3-1)} = \frac{N}{3} = \frac{10}{3}$$

$$E(T_4) = \frac{2N}{4(4-1)} = \frac{N}{6} = \frac{10}{6}$$

present

The expected **waiting time** to go from i lineages to 1 lineage (TMRCA) is the sum overall all coalescent intervals

past



$$E(\text{TMRCA}) =$$

$$E(T_4) + E(T_3) + E(T_2) =$$

$$\frac{N}{6} + \frac{N}{3} + N =$$

$$\frac{3N}{2} = \frac{30}{2}$$

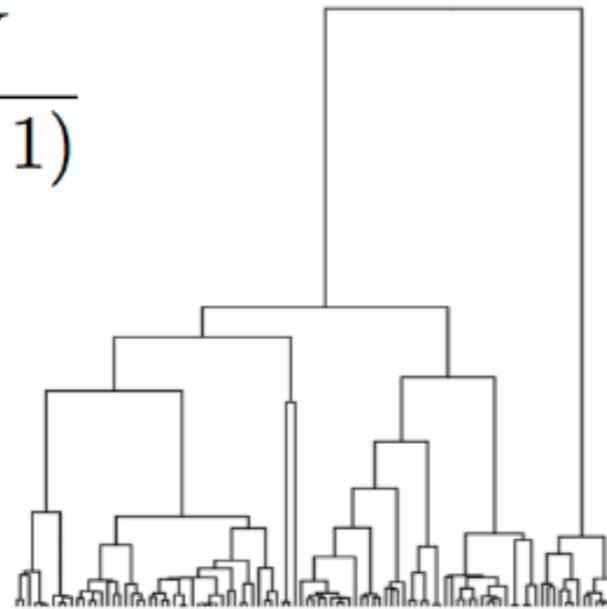
= 15 generations

The expected **waiting time** to go from i lineages to 1 lineage is the sum overall all coalescent intervals:

$$T_{\text{MRCA}} = \sum_{i=2}^n T_i \quad E[T_i] = \frac{2N}{i(i-1)}$$

$$\begin{aligned} E[T_{\text{MRCA}}] &= \sum_{i=2}^n \frac{2N}{i(i-1)} \\ &= 2N \sum_{i=2}^n \left(\frac{1}{i-1} - \frac{1}{i} \right) \\ &= 2N \left(1 - \frac{1}{2} + \frac{1}{2} - \cdots - \frac{1}{n-1} + \frac{1}{n-1} - \frac{1}{n} \right) \\ &= 2N \left(1 - \frac{1}{n} \right) \end{aligned}$$

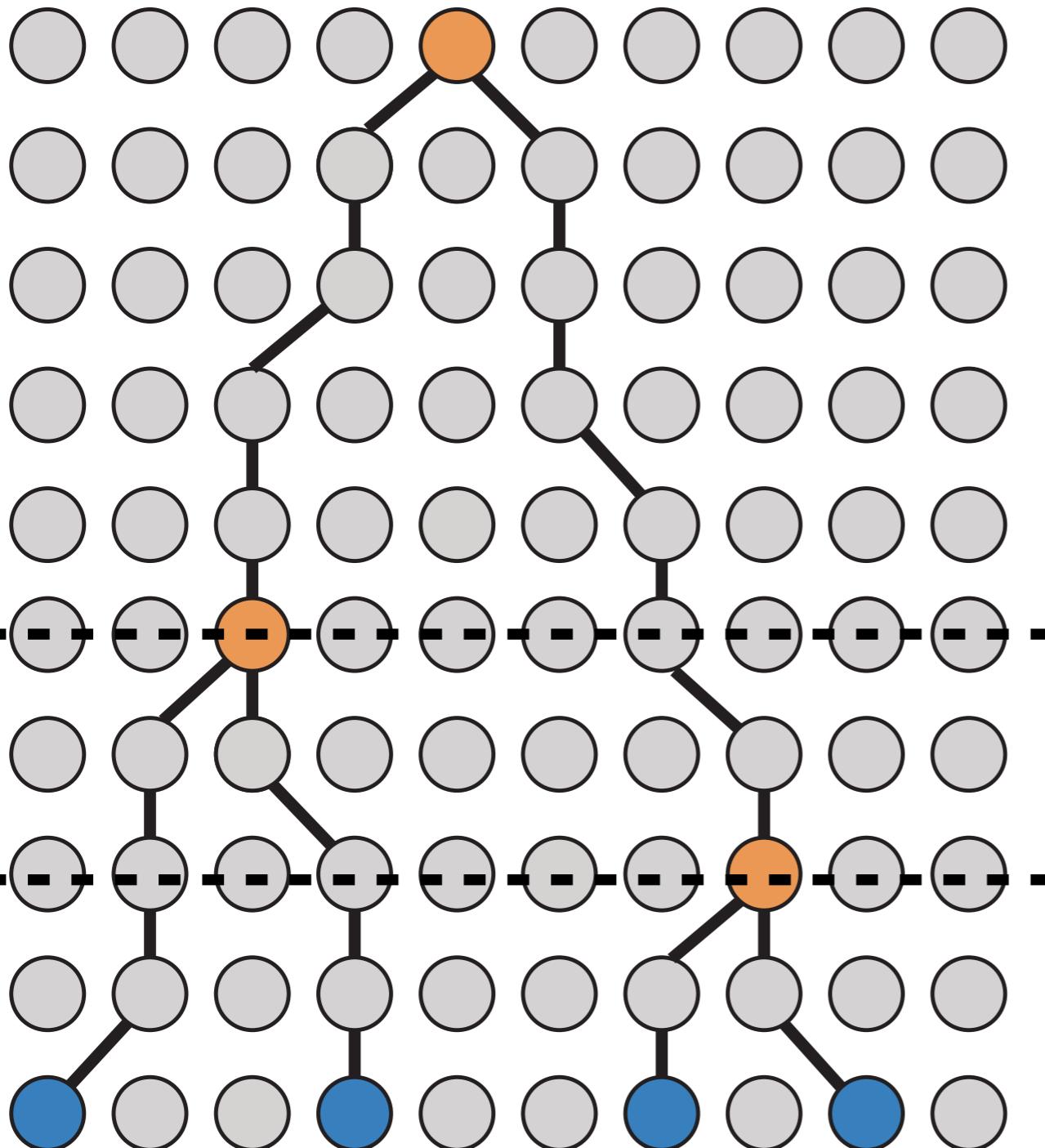
$$E[T_{\text{MRCA}}] = \lim_{n \rightarrow \infty} 2N$$



* Bedford SISMID
lecture, Coalescent 57

The expected TMRCA approaches 2N

past



N (population size) = 10
n (sampled individuals) = 4

$$E(\text{TMRCA}) = 2N \left(1 - \frac{1}{n}\right)$$

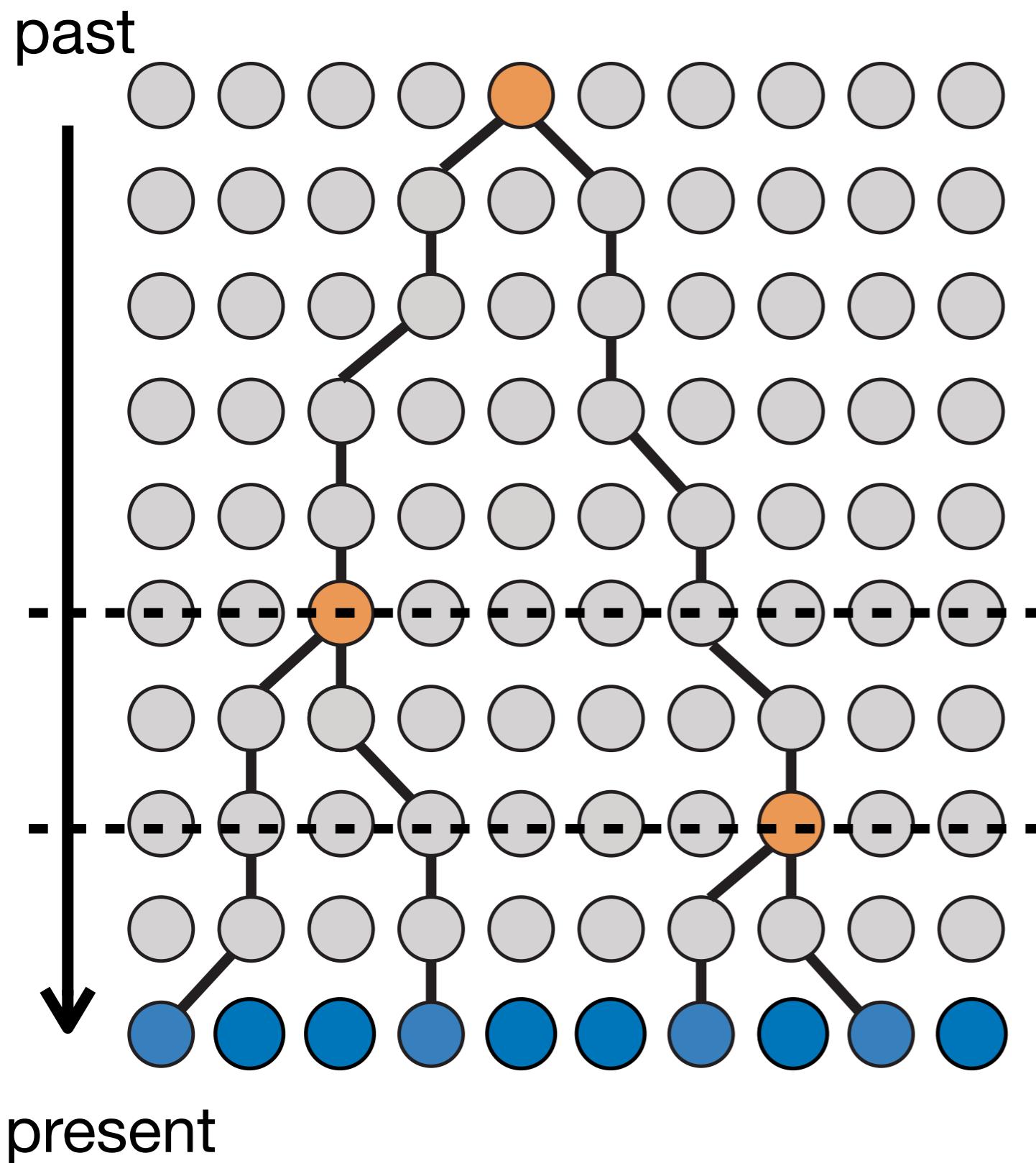
$$= 2 * 10 \left(1 - \frac{1}{4}\right)$$

$$= 20(0.75)$$

= **15 generations**

present

The expected TMRCA approaches 2N



N (population size) = 10
n (sampled individuals) = 10

$$\begin{aligned} E(\text{TMRCA}) &= 2N \left(1 - \frac{1}{n}\right) \\ &= 2 * 10 \left(1 - \frac{1}{10}\right) \\ &= 20(0.90) \\ &= \mathbf{18 \text{ generations}} \end{aligned}$$

A few key points about coalescent trees:

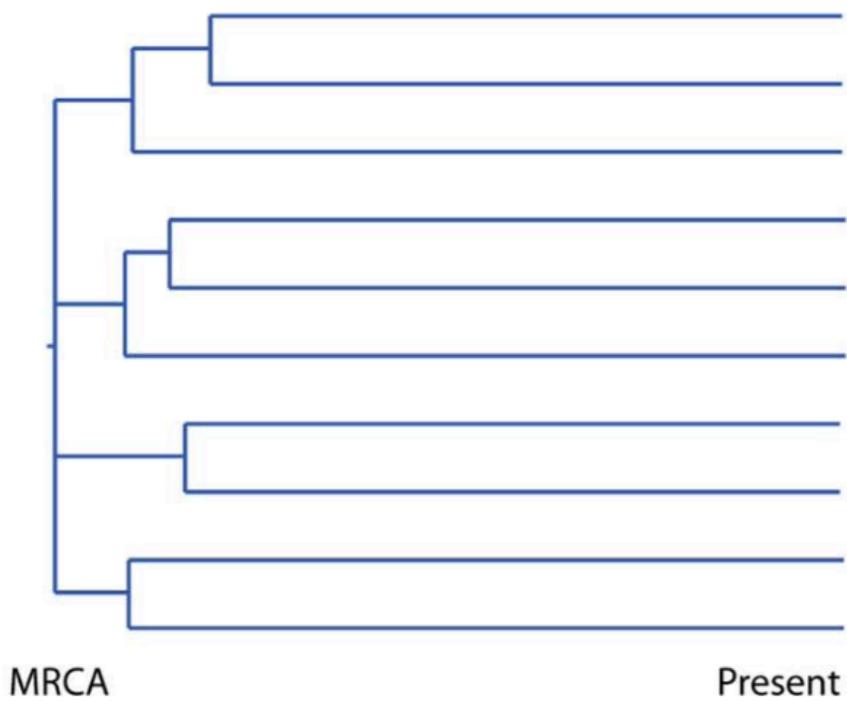
1. Coalescence takes longer in large populations
2. Coalescence takes longer when there are fewer lineages present
3. Most of the tree height is contributed by only 2 lineages
4. The TMRCA for a set of sampled lineages is $\sim 2N$
5. The coalescent gives us a framework for setting an expectation for genealogies under different evolutionary scenarios

The coalescent is conceived under a Wright-Fisher model, which makes a lot of assumptions

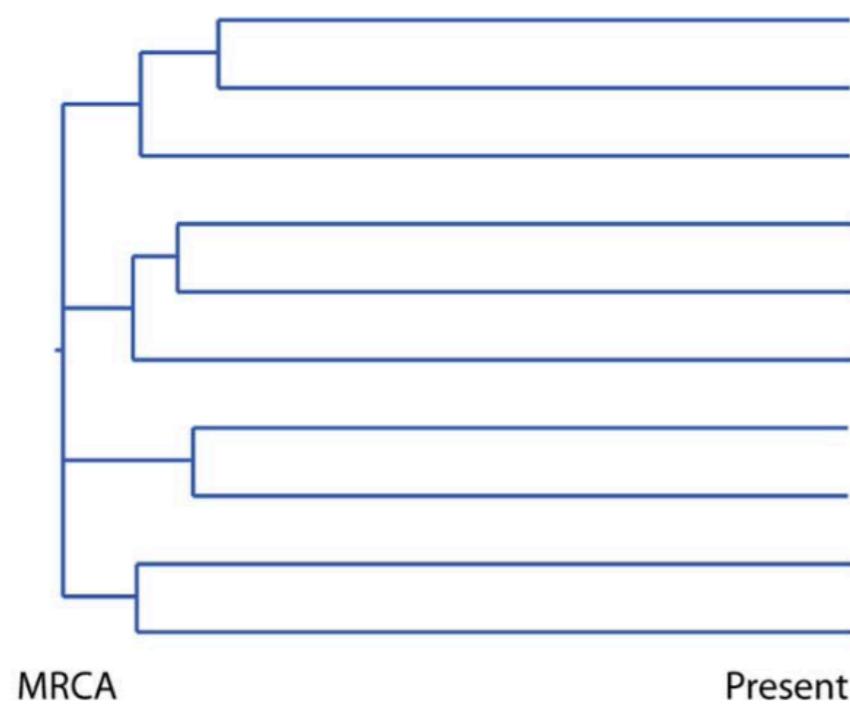
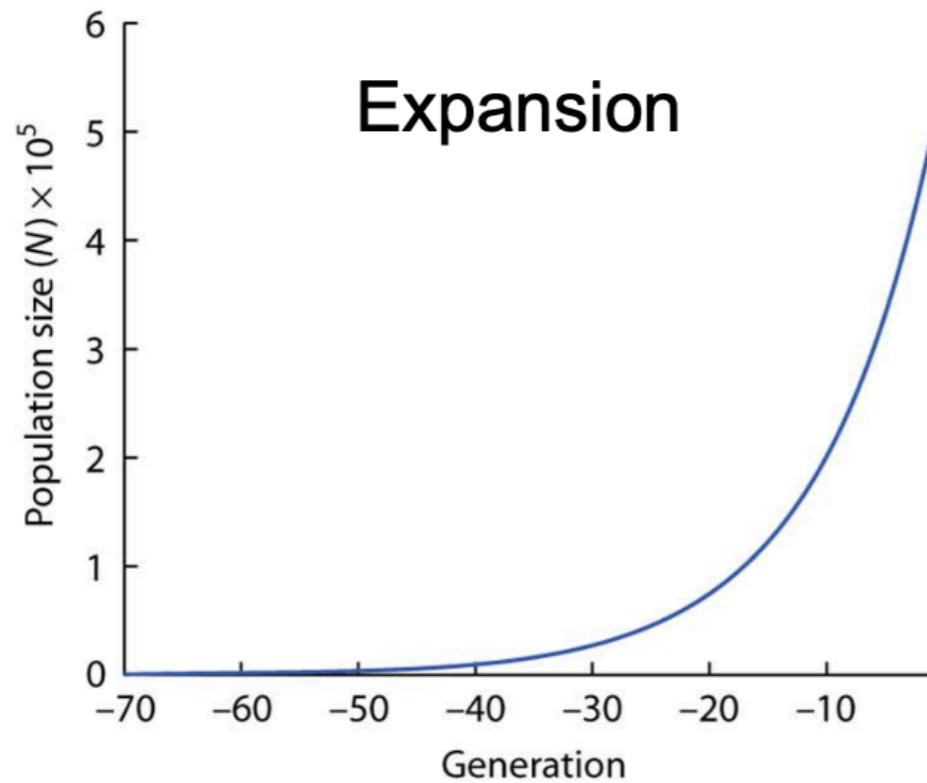
- No selection, migration, or recombination
- Random mating
- Population size is constant over time
- Sampled lineages $\ll N$; maximum of 1 coalescent event per generation

We can use our knowledge of expected coalescence times to find evidence of deviation from these assumptions!

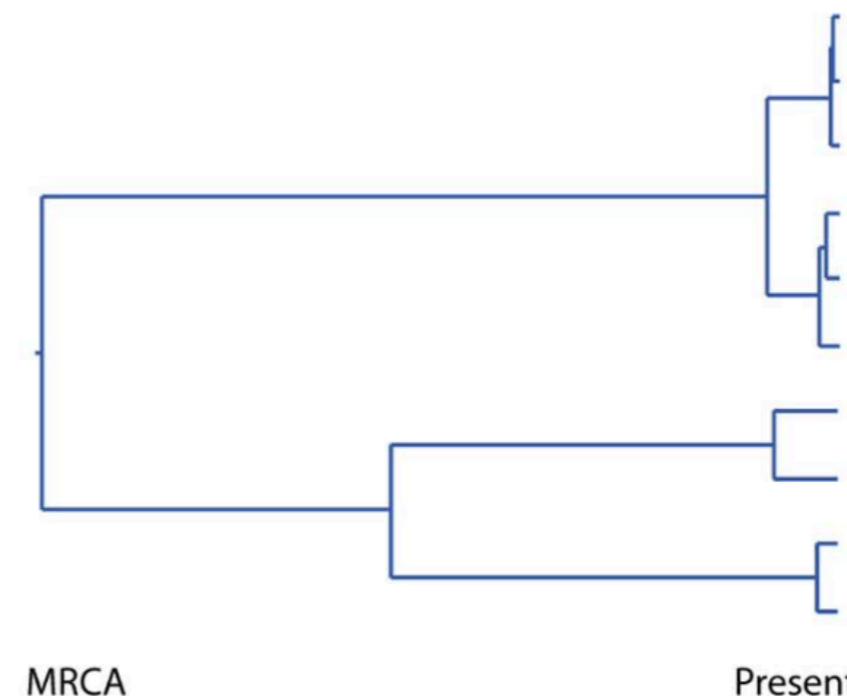
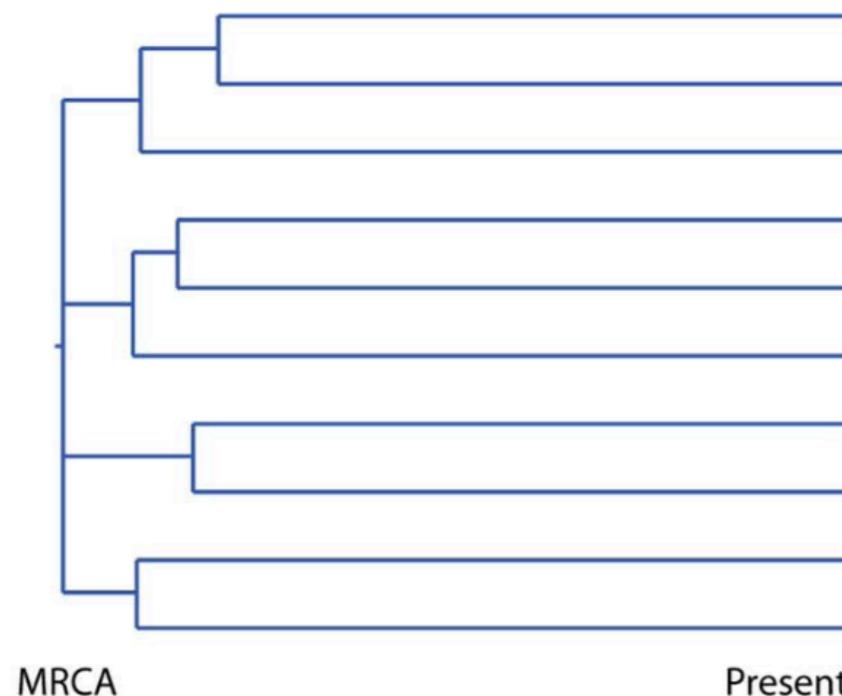
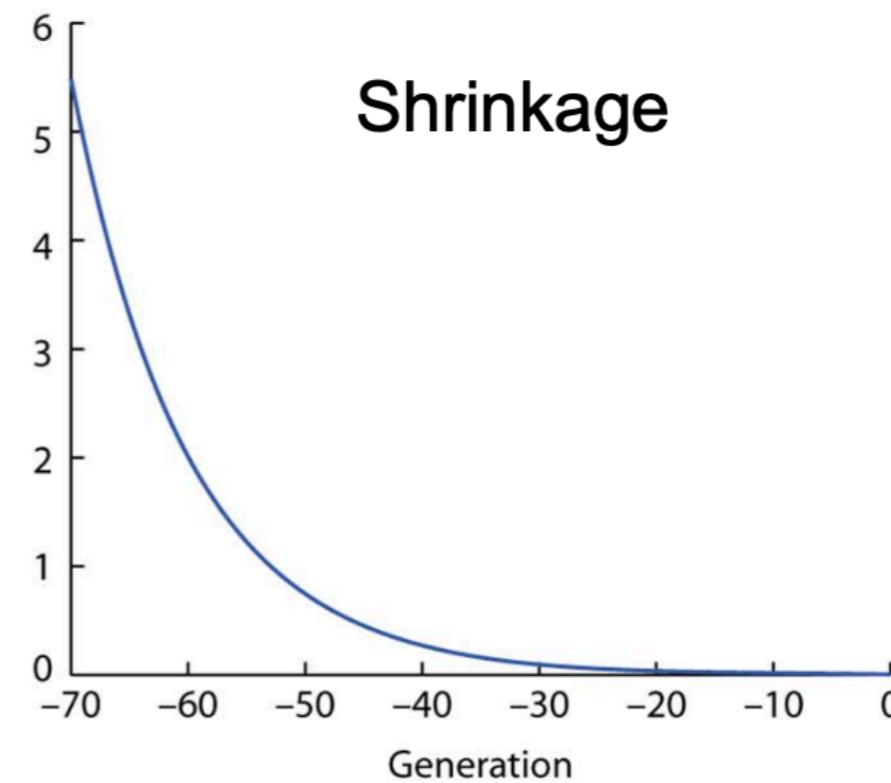
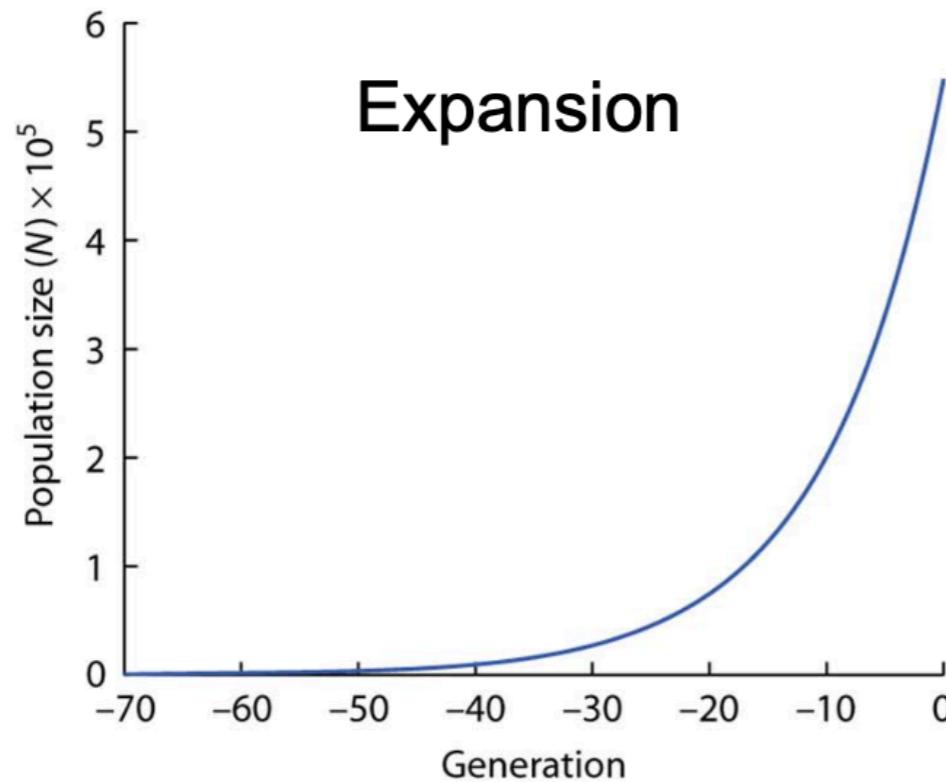
Coalescent trees look very different in large and small populations

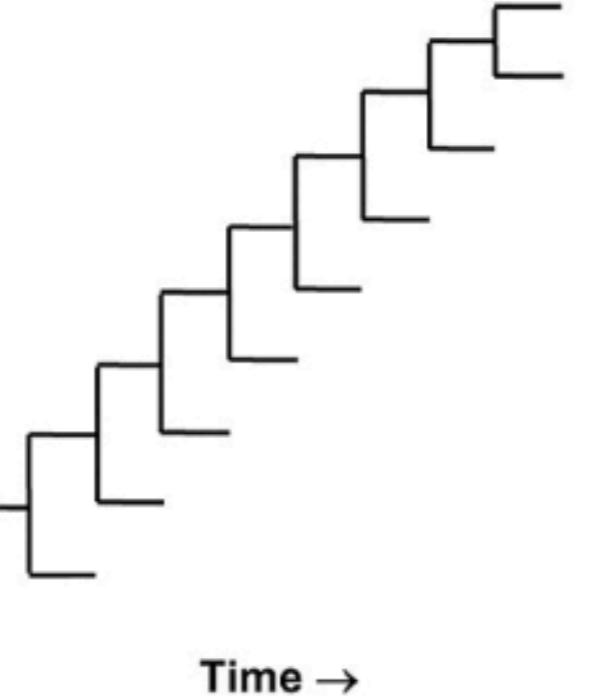
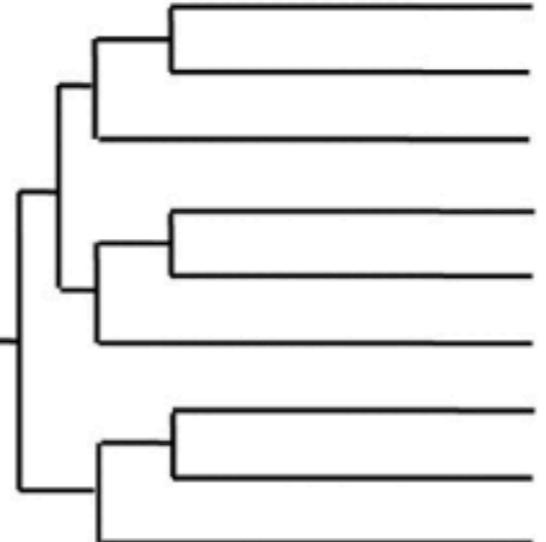
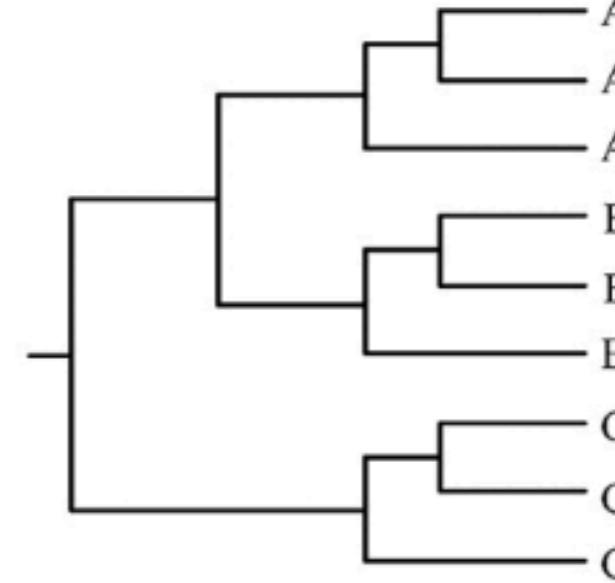
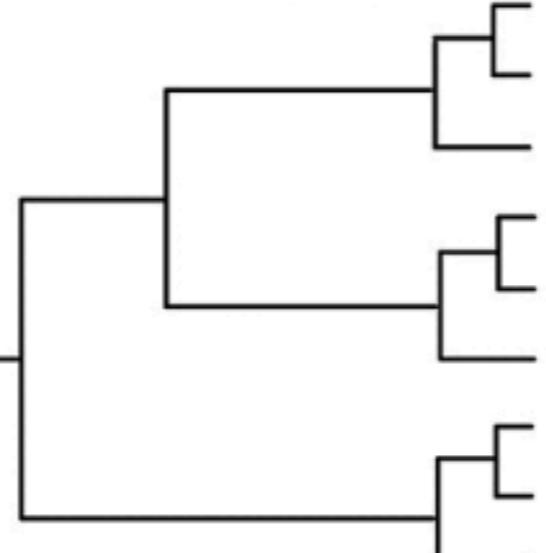
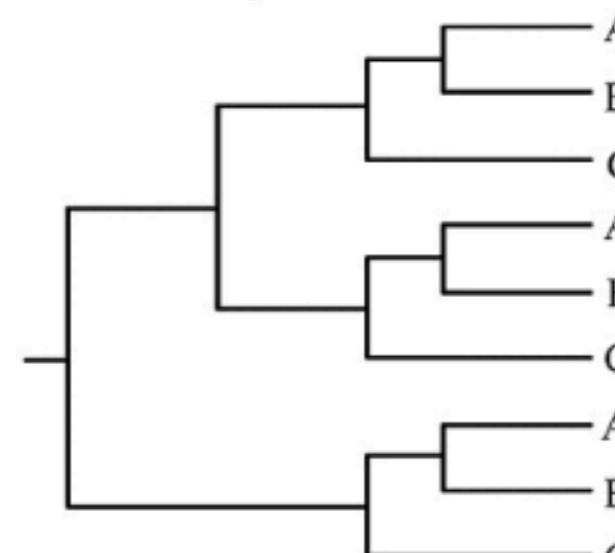


Coalescent trees look very different in large and small populations



Coalescent trees look very different in large and small populations



Continual Immune Selection		Weak or Absent Immune Selection	
		Tree shape controlled by non-selective population dynamic processes	
		Population size dynamics	Spatial dynamics
Idealized Phylogeny Shapes		Exponential growth 	Strong spatial structure 
		Constant size 	Weak spatial structure 
Examples	Human influenza A virus intra-host HIV	inter-host HIV inter-host HCV	Measles, rabies inter-host HIV
Tree Inferences	Detection of antigenic escape mutations	Estimation of population growth rates	Estimation of population migration rates

The key points we'll be making today:

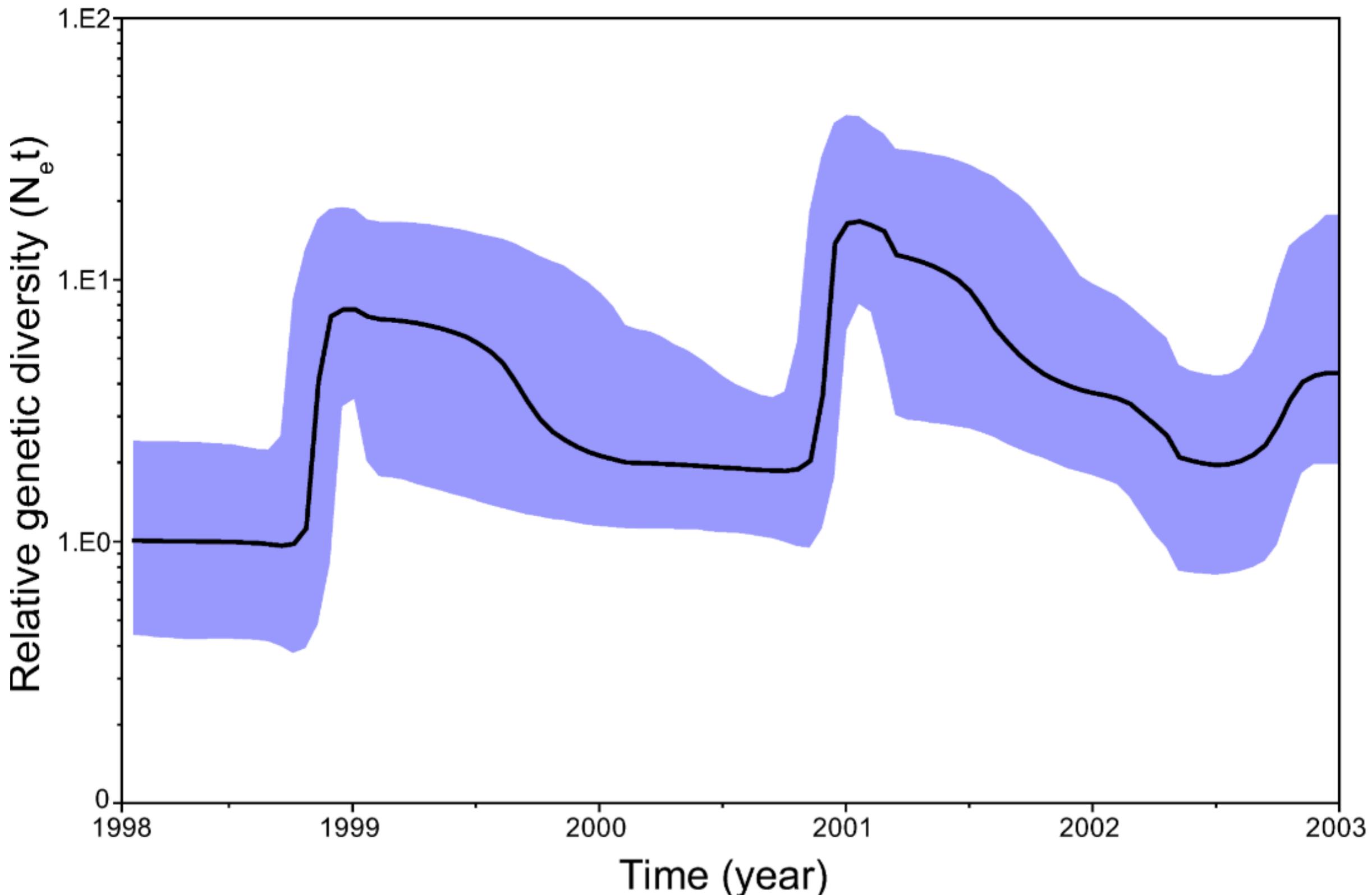
1. Pathogen evolution happens on the same timescale as transmission, meaning that for rapidly evolving pathogens, **evolution, ecology, and epidemiology are linked**.
2. **Coalescent theory** allows us to generate expectations for phylogenies generated by populations that are evolving in particular ways.
3. These features allow us to use trees to infer parameters, and vice versa. We will see **examples** of how others have done this.
4. How can **we apply** these concepts to our own questions?

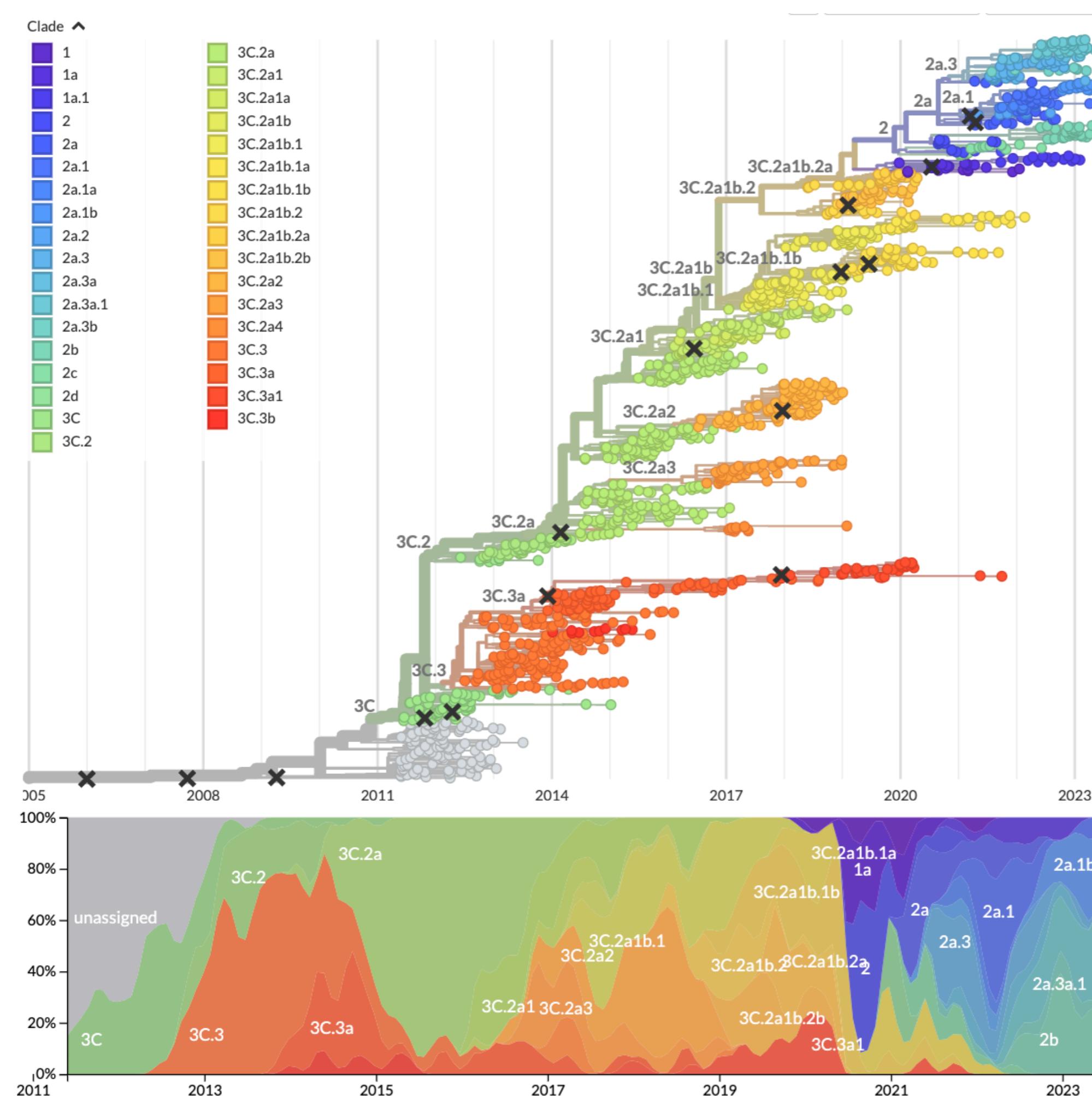
The key points we'll be making today:

1. Pathogen evolution happens on the same timescale as transmission, meaning that for rapidly evolving pathogens, **evolution, ecology, and epidemiology are linked**.
2. **Coalescent theory** allows us to generate expectations for phylogenies generated by populations that are evolving in particular ways.
3. These features allow us to use trees to infer parameters, and vice versa. We will see **examples** of how others have done this.
4. How can **we apply** these concepts to our own questions?

Some real life examples

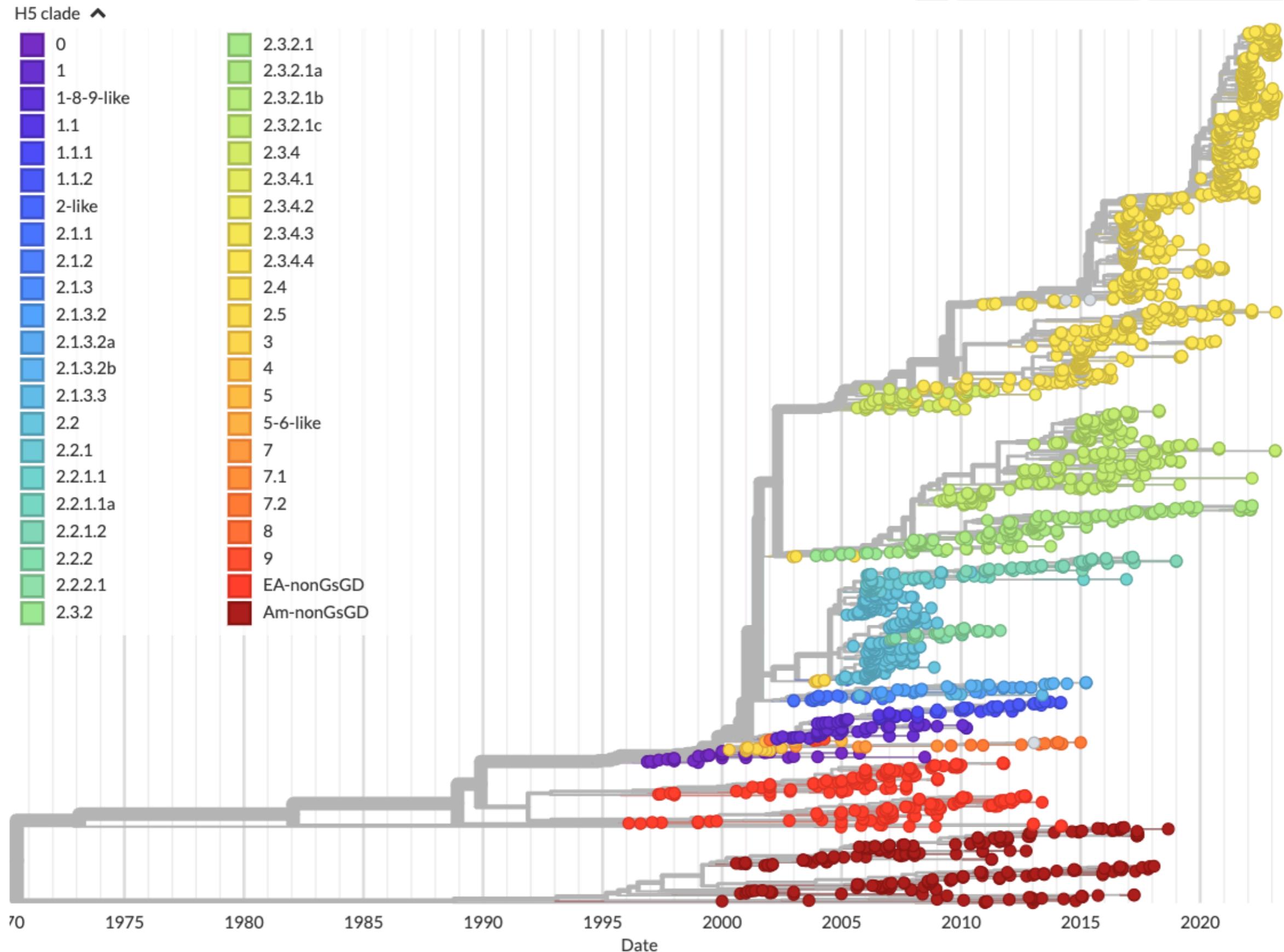
Population size changes in influenza in New York

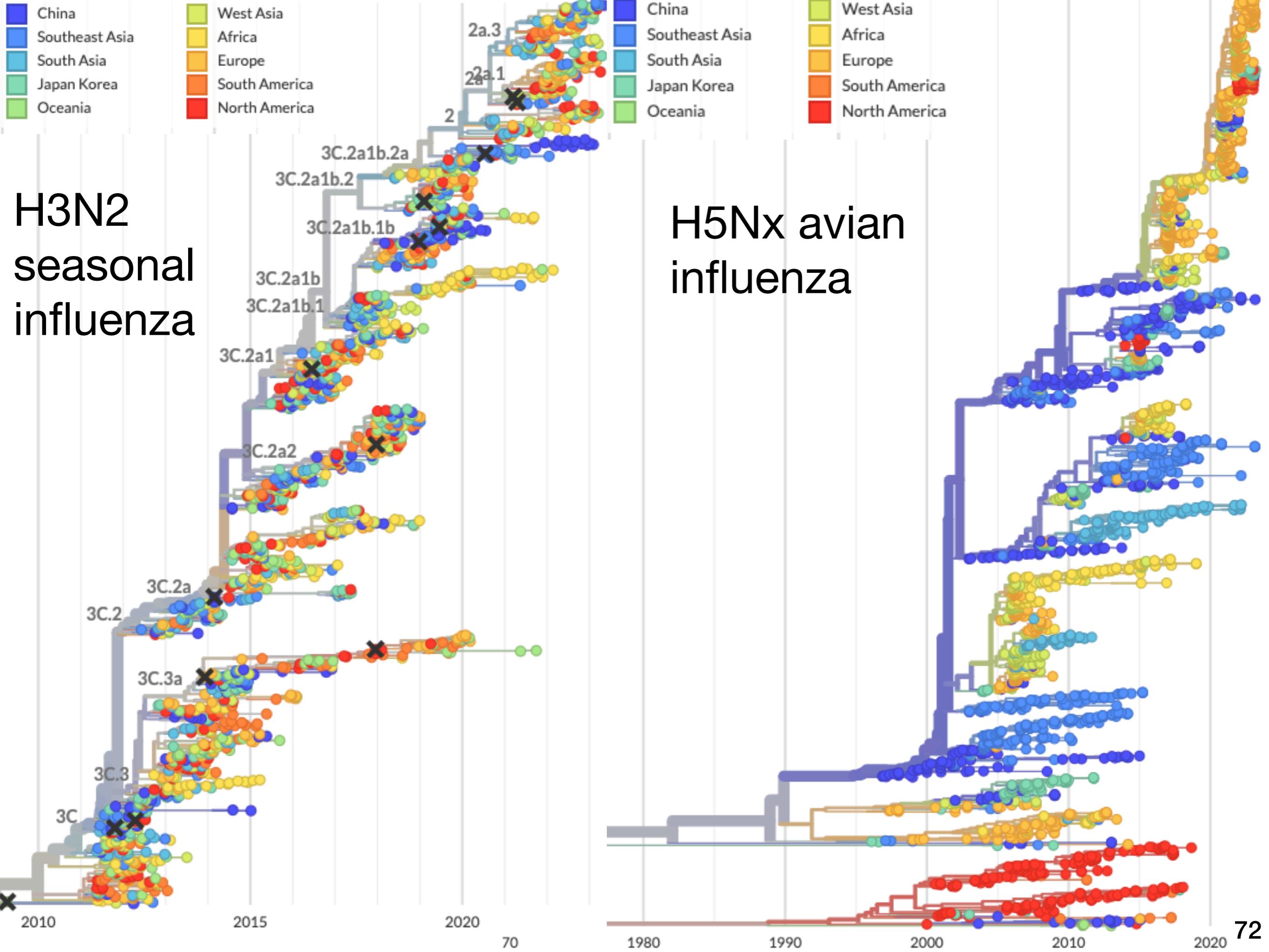




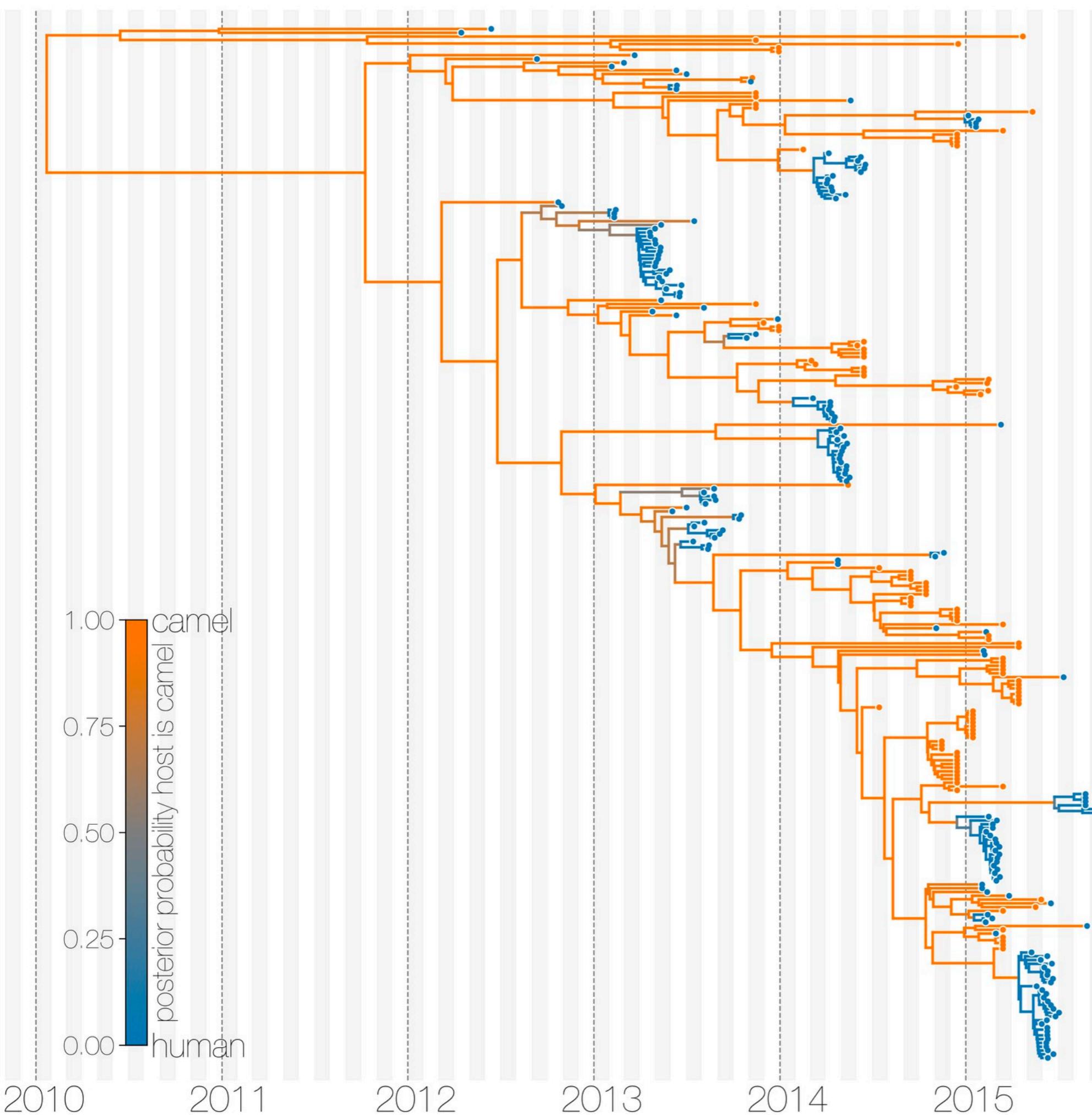
Seasonal influenza evolves yearly, necessitating vaccine updates

Avian influenza evolves with less selection, and more geographic structure





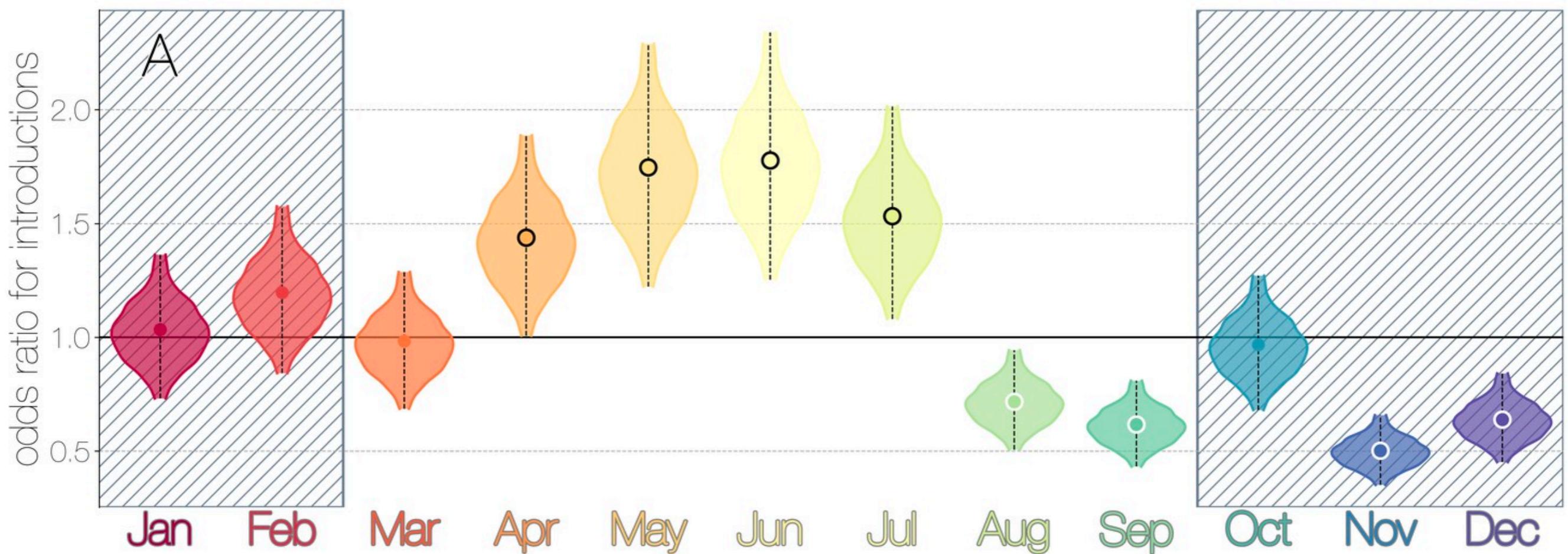
Inferring host switching for MERS



NeTau estimates:
Humans:
0.24 years
Camels: 3.5 years

Dudas et al, eLife
2020

Spring calving season correlates with an increase in the number of cross-species transmission events



Extra resources and cool examples:

- West Nile virus in the Americas: <https://nextstrain.org/narratives/twenty-years-of-WNV> and <https://journals.plos.org/plospathogens/article?id=10.1371/journal.ppat.1008042>
- Ebola in West Africa: <https://www.nature.com/articles/nature22040>
- HIV in humans: <https://www.science.org/doi/full/10.1126/science.1256739>
- Cryptic transmission of SARS-CoV-2 in Washington state: <https://www.science.org/doi/full/10.1126/science.abc0523>
-

We will be using 2 pieces of software to explore phylodynamic inference



HELP DOCS BLOG LOGIN

Nextstrain

Real-time tracking of pathogen evolution

Nextstrain is an open-source project to harness the scientific and public health potential of pathogen genome data. We provide a continually-updated view of publicly available data alongside powerful analytic and visualization tools for use by the community. Our goal is to aid epidemiological understanding and improve outbreak response. If you have any questions, or simply want to say hi, please give us a shout at hello@nextstrain.org.

[READ MORE](#)



How can we scale this to more than 2 lineages?

In any given generation, the **probability** that i individuals pick the same parent and coalesce:

$$\binom{i}{2} \left(\frac{1}{N}\right)$$

The expected **waiting time** i lineages to coalesce is exponentially distributed with a mean of:

$$\frac{1}{\binom{i}{2}} (N)$$

What the heck is i choose 2?

In any given generation, the **probability** that i individuals pick the same parent and coalesce:

$$\left(\frac{i(i-1)}{2} \right) \left(\frac{1}{N} \right)$$

The expected **waiting time** i lineages to coalesce is exponentially distributed with a mean of:

$$\frac{1}{\left(\frac{i(i-1)}{2} \right) (N)}$$