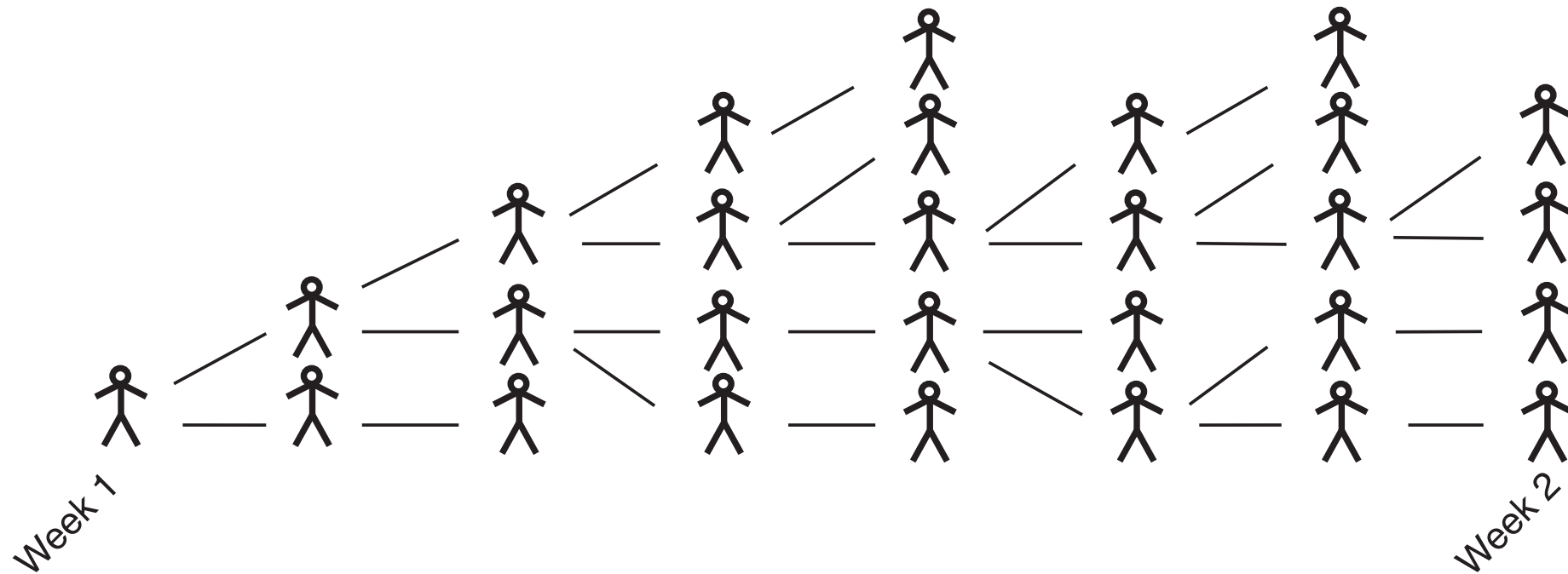


Accounting for population structure

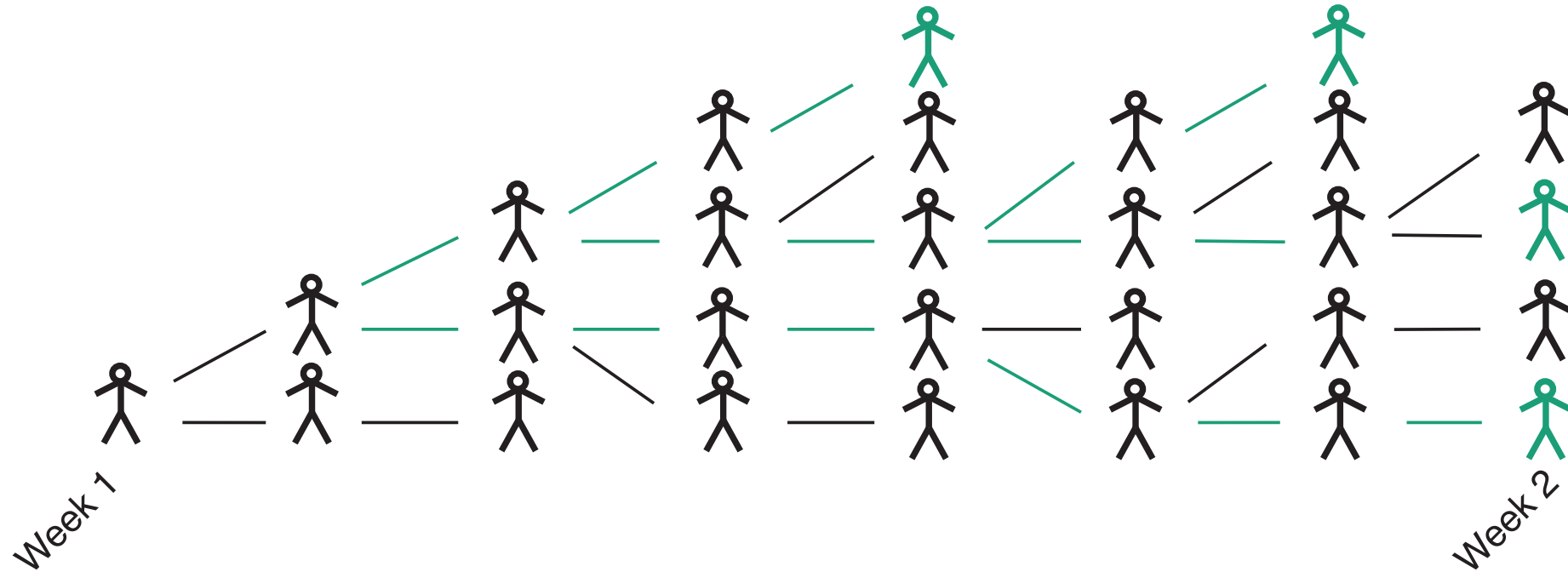
Nicola F. Müller

e-mail: nicola.mueller@ucsf.edu

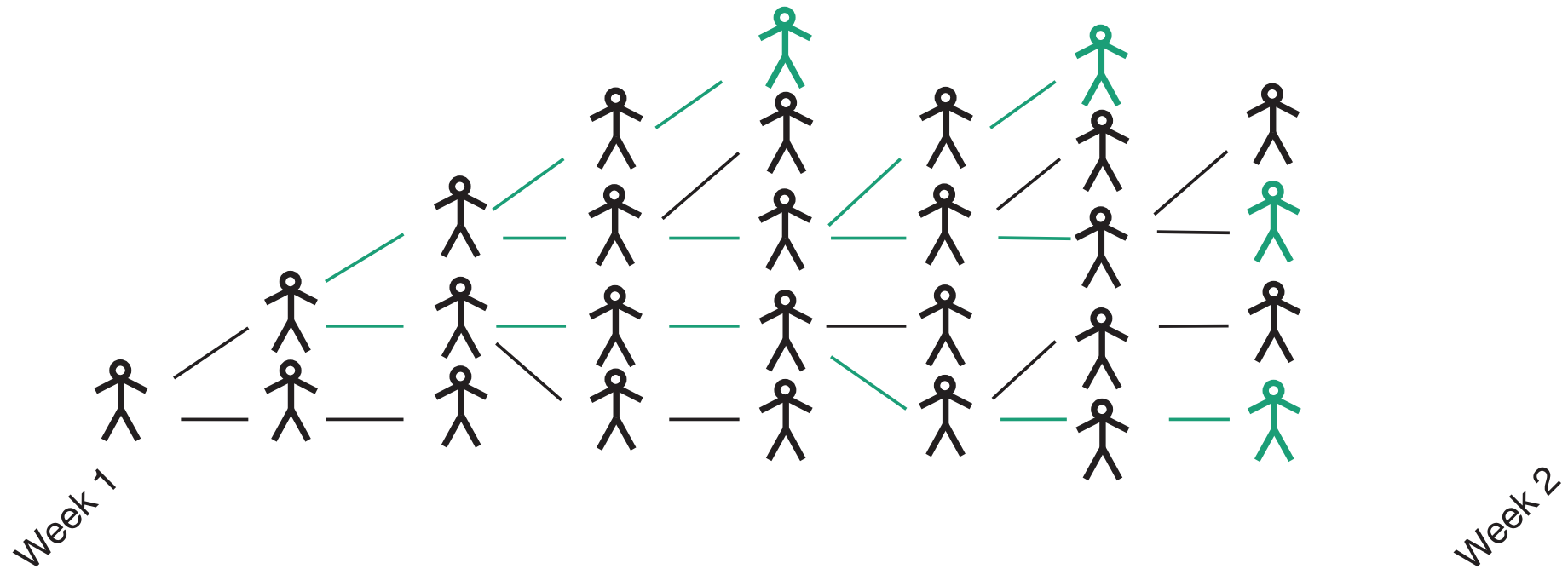
Genomics allows tracking the shared ancestry.



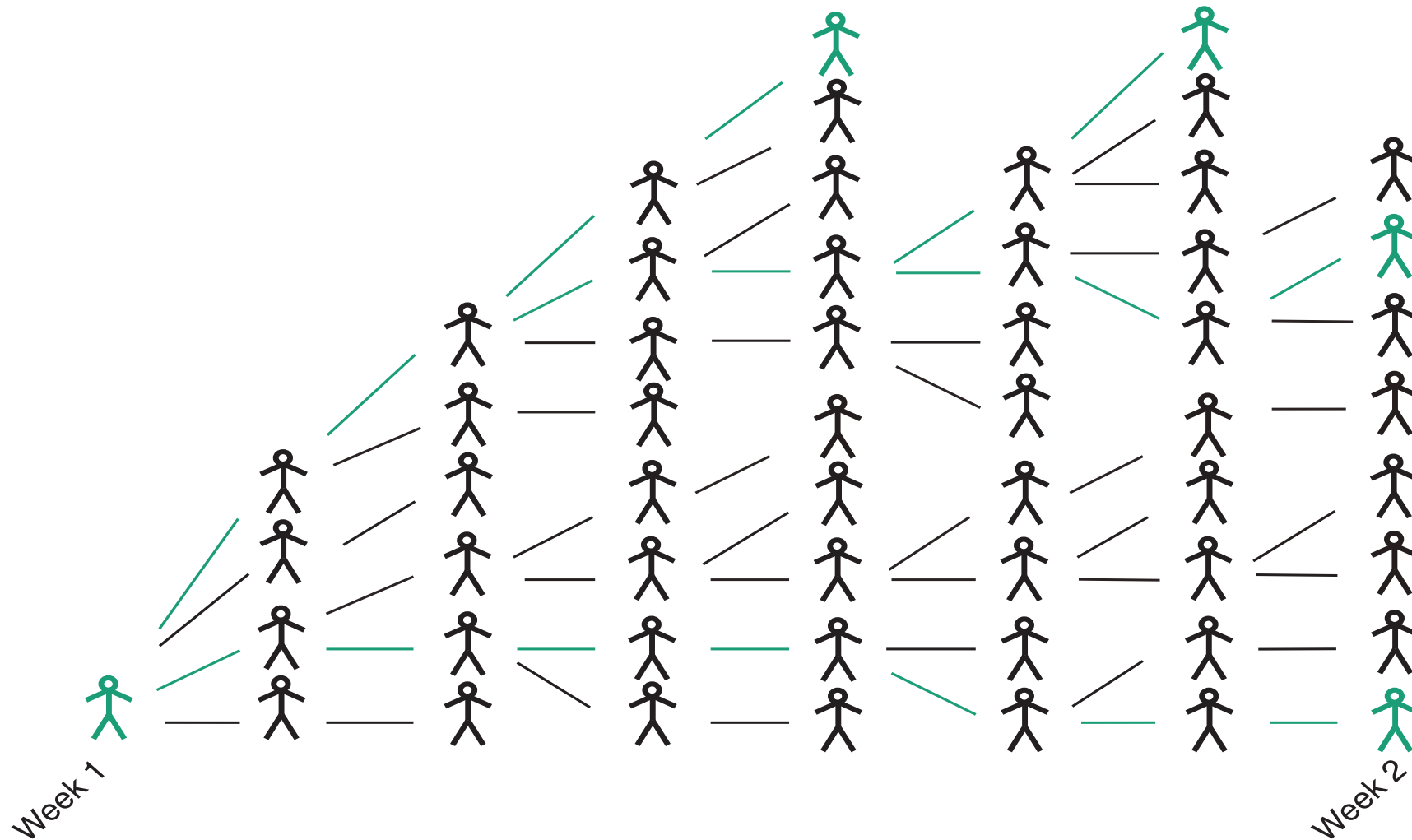
If we track the shared ancestry of the 4 sampled individuals, we will reach a common ancestor more quickly, the smaller the N_e .



What happens if we have the same number of infected individuals, but higher turnover?



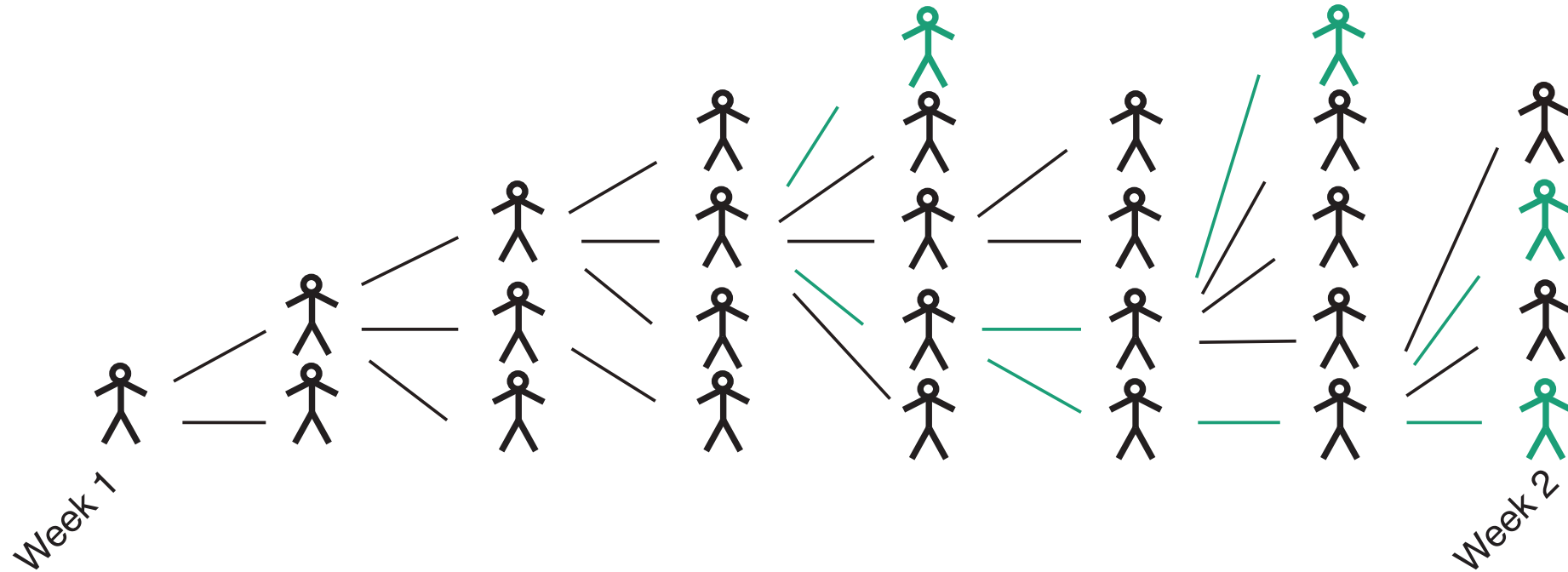
What happens if we have a higher number of infected individuals and higher turnover?



The effective population size (N_e) can be denoted as a function of epi-parameters for SIR models

$$N_e(t) = \frac{I(t)}{\theta \frac{S(t)}{N}}$$

What happens if we have a skewed offspring distribution?

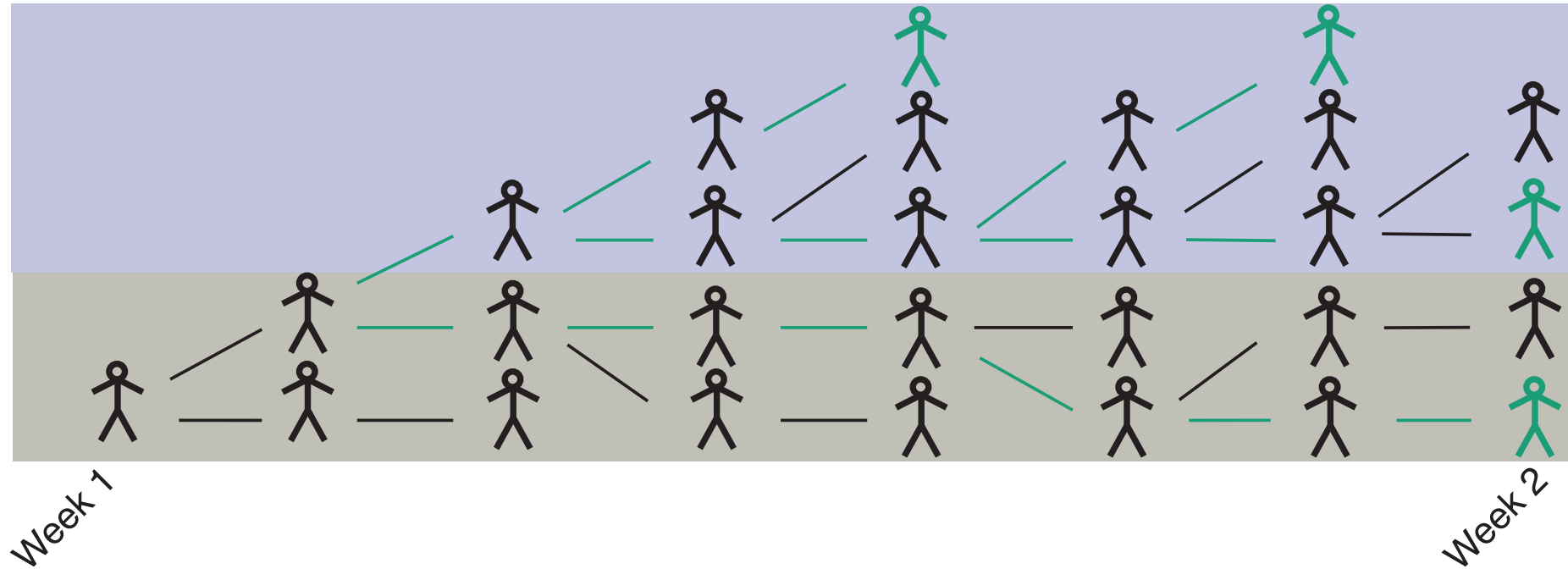


We can also account for some offspring distributions in the Ne

$$N_e(t) = \frac{I(t)}{\theta \left(1 + \frac{1}{k}\right)}$$

Li, L. M., Grassly, N. C., & Fraser, C. (2017). *“Quantifying Transmission Heterogeneity Using Both Pathogen Phylogenies and Incidence Time Series.”* Molecular Biology and Evolution, 34(11), 2982-2995. <https://doi.org/10.1093/molbev/msx195>

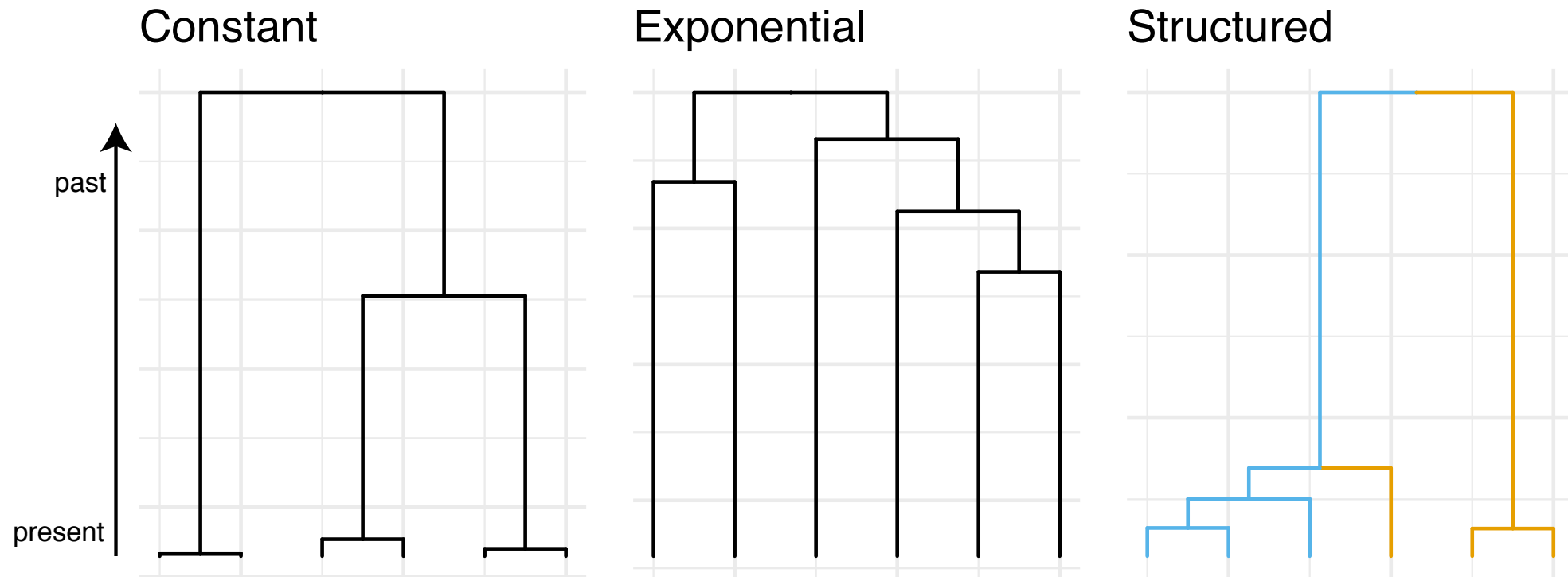
What happens if there is population structure?



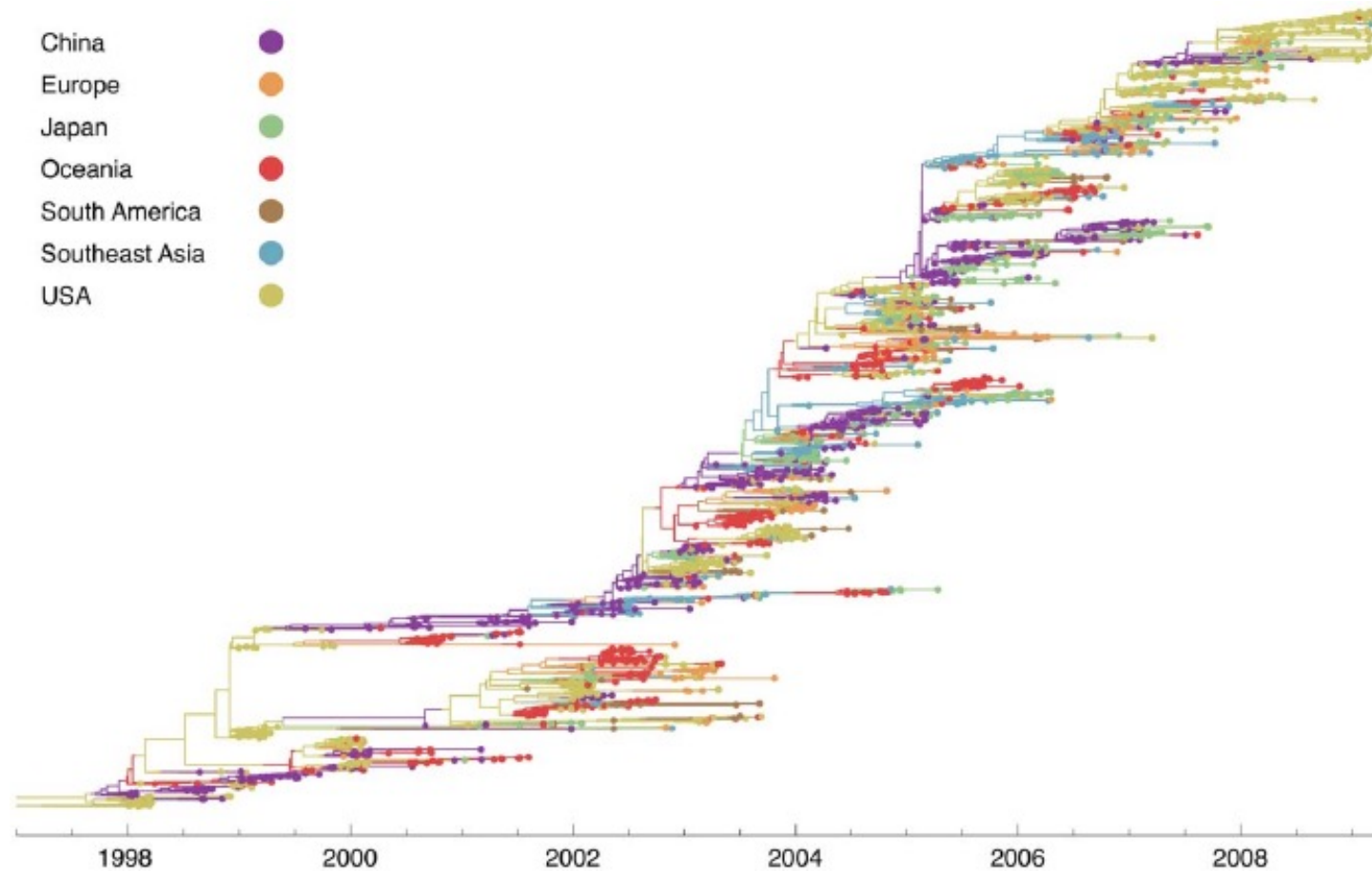
Once there is population structure, the meaning of the effective population size is not obvious

$$N_e(t) = ?$$

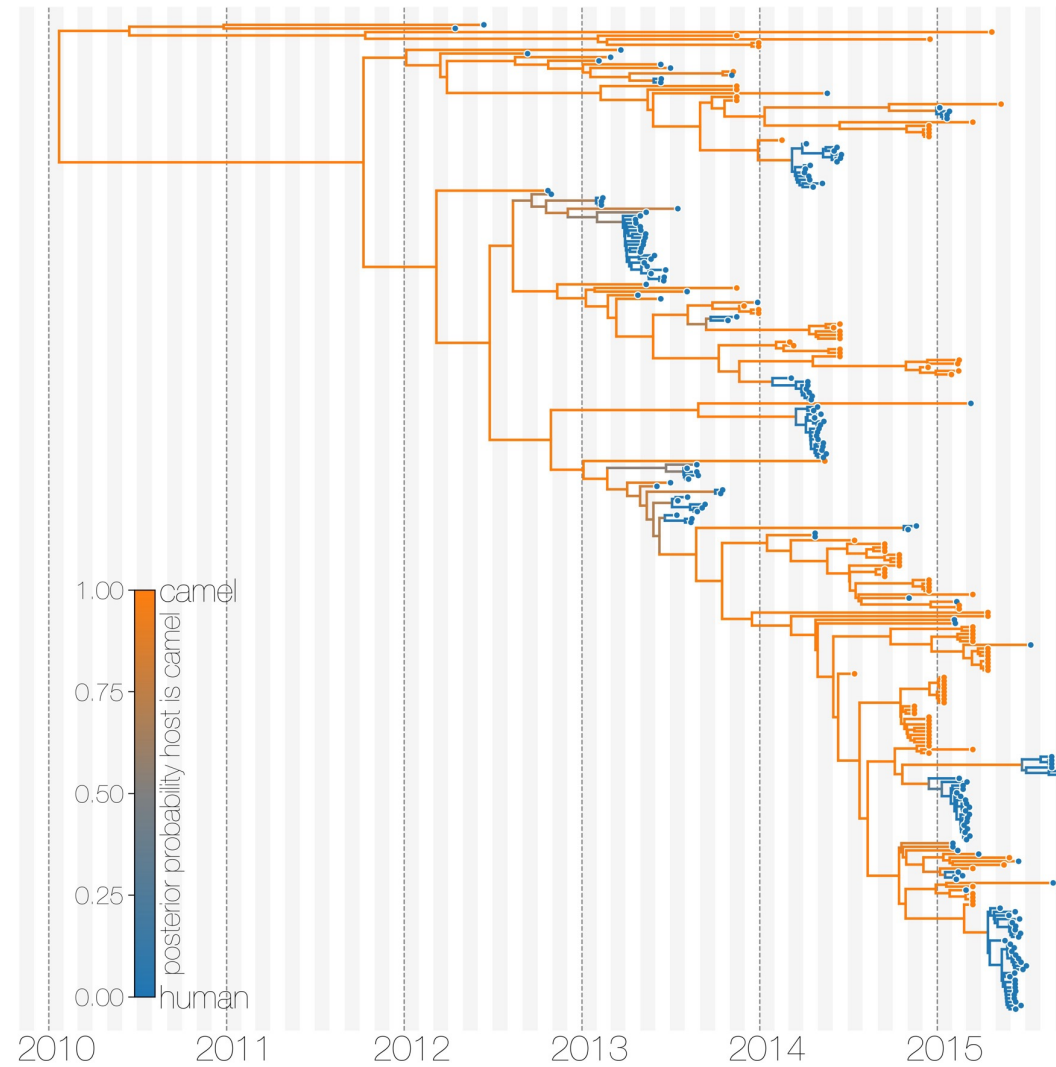
Phylogenetic trees are formed by population processes and contain information about them



H3N2 sequences show spatial clustering



MERS sequences cluster by host of isolation



In the posterior probability, modelling the population dynamics happens in the tree prior.

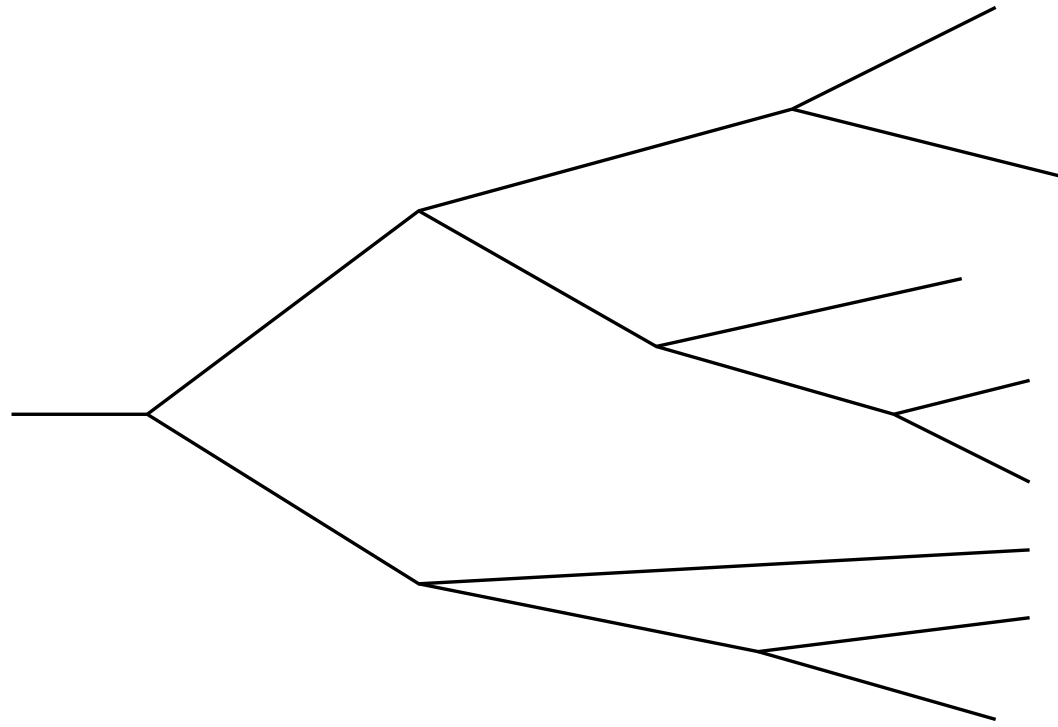
$$P(\text{Tree} \mid \text{Data}) = \frac{P(\text{Data} \mid \text{Tree}) P(\text{Tree}) P(\text{Model}) P(\text{Prior}) P(\text{Clock})}{P(\text{Data})}$$

The equation represents the posterior probability of the tree given the data. The numerator consists of five terms: the likelihood of the data given the tree, the prior probability of the tree, the probability of the model, the prior probability of the model, and the probability of the clock. The denominator is the marginal likelihood of the data.

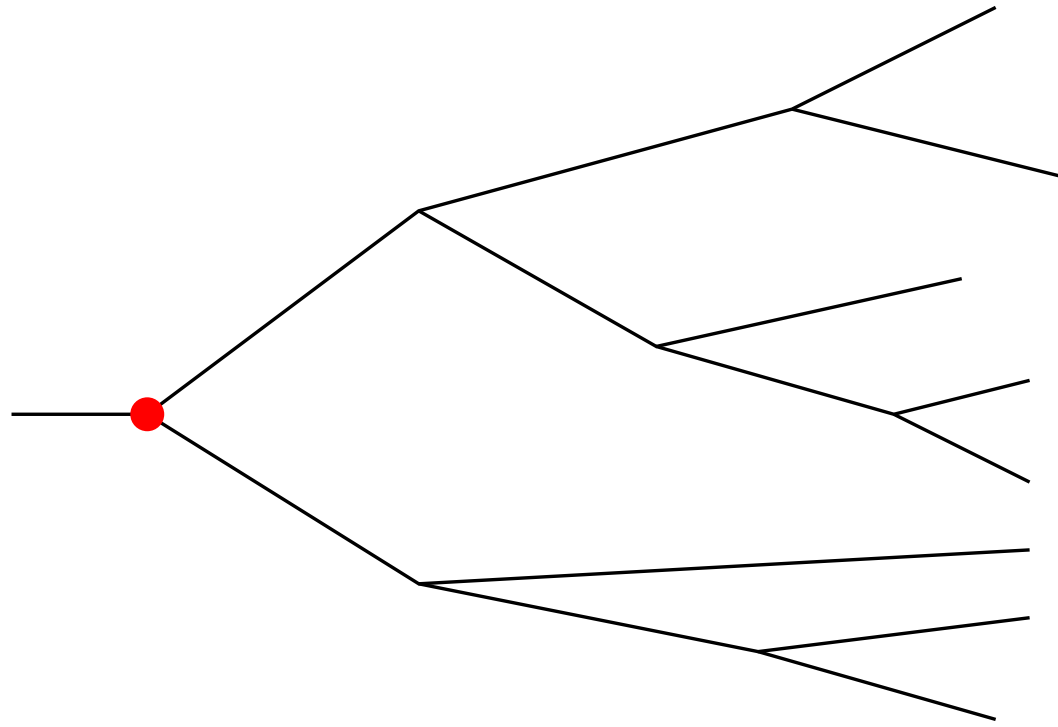
Models to account for population structure:

- **DTA:** Assumes structure to be a neutral trait evolving “on top” of a phylogenetic tree (Lemey 2009).
- **Structured Coalescent:** Models how lineages coalesce within and migrate between discrete locations.
- **Multi-type Birth-death models:** Models how lineages give birth to other lineage in the same and migrate between them.

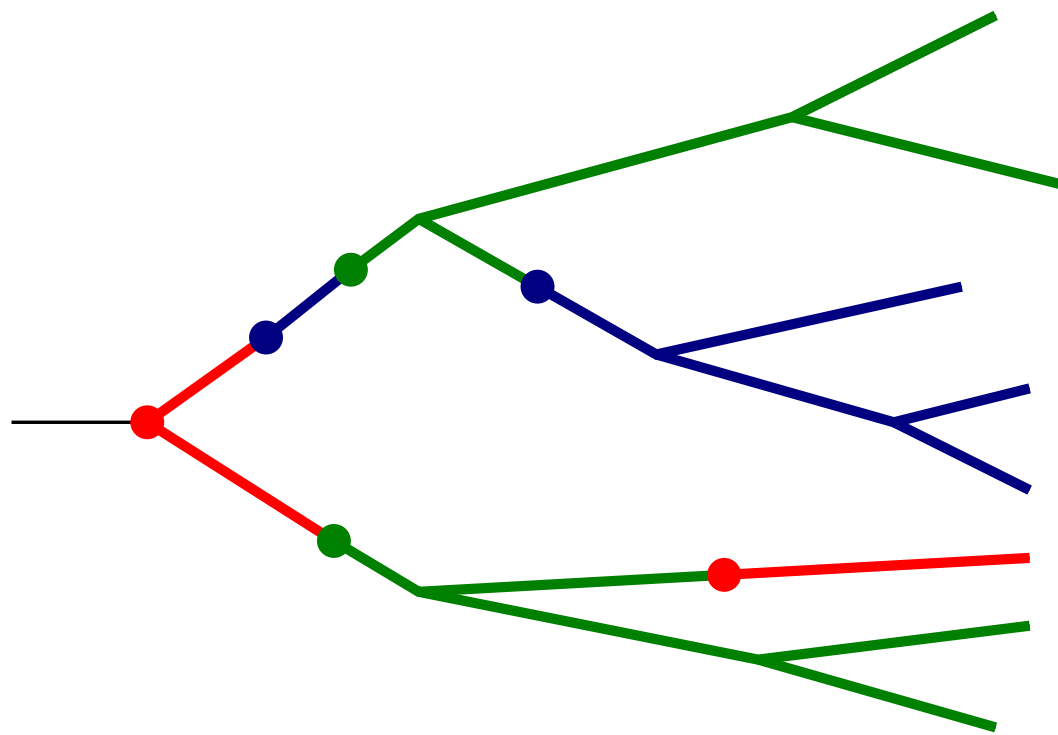
Tree was generated by a different process (e.g. a coalescent process).



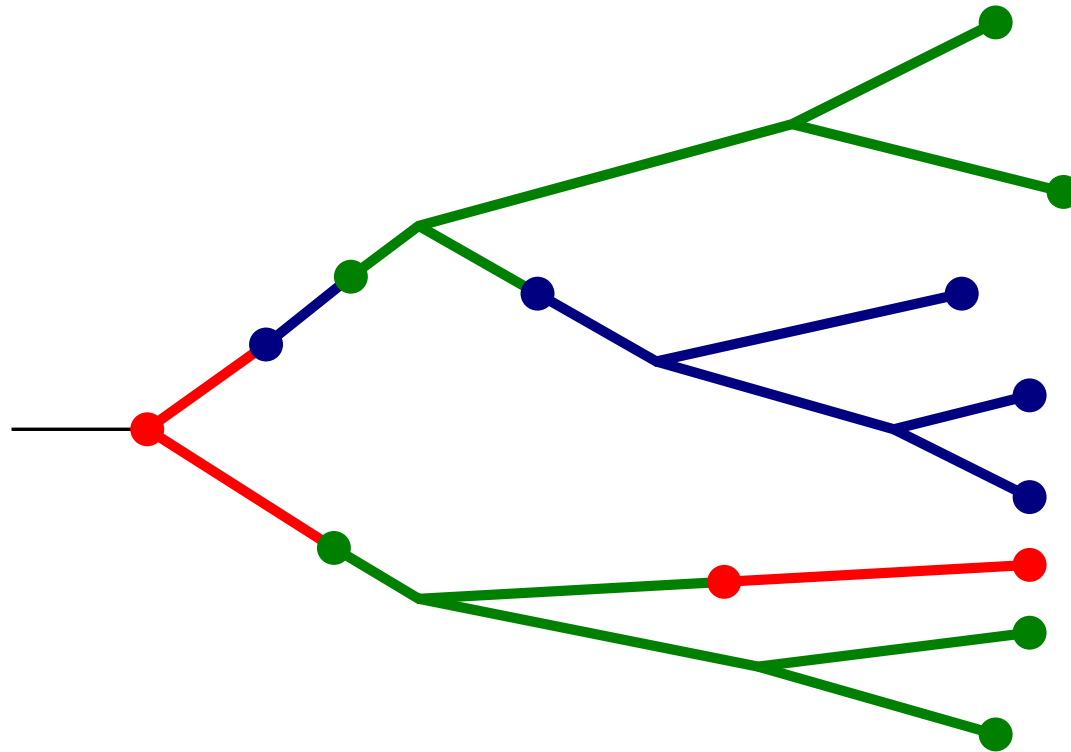
Root of the tree is in some location (here red)



Then models the migration process on top of this tree
from the root to the tips



Sampling locations are considered Data.



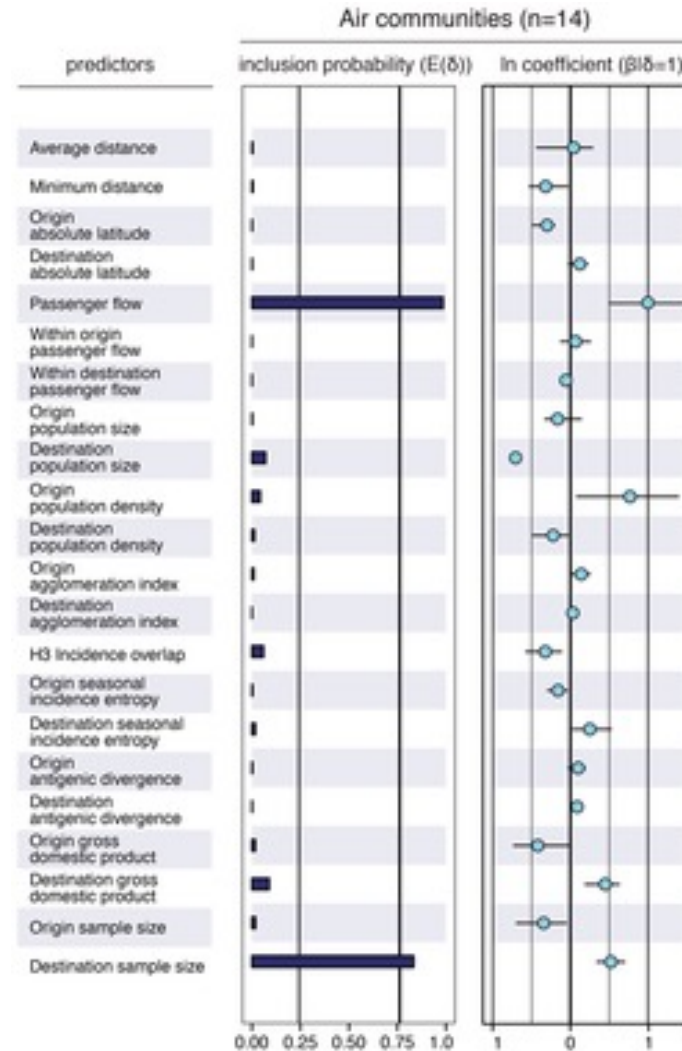
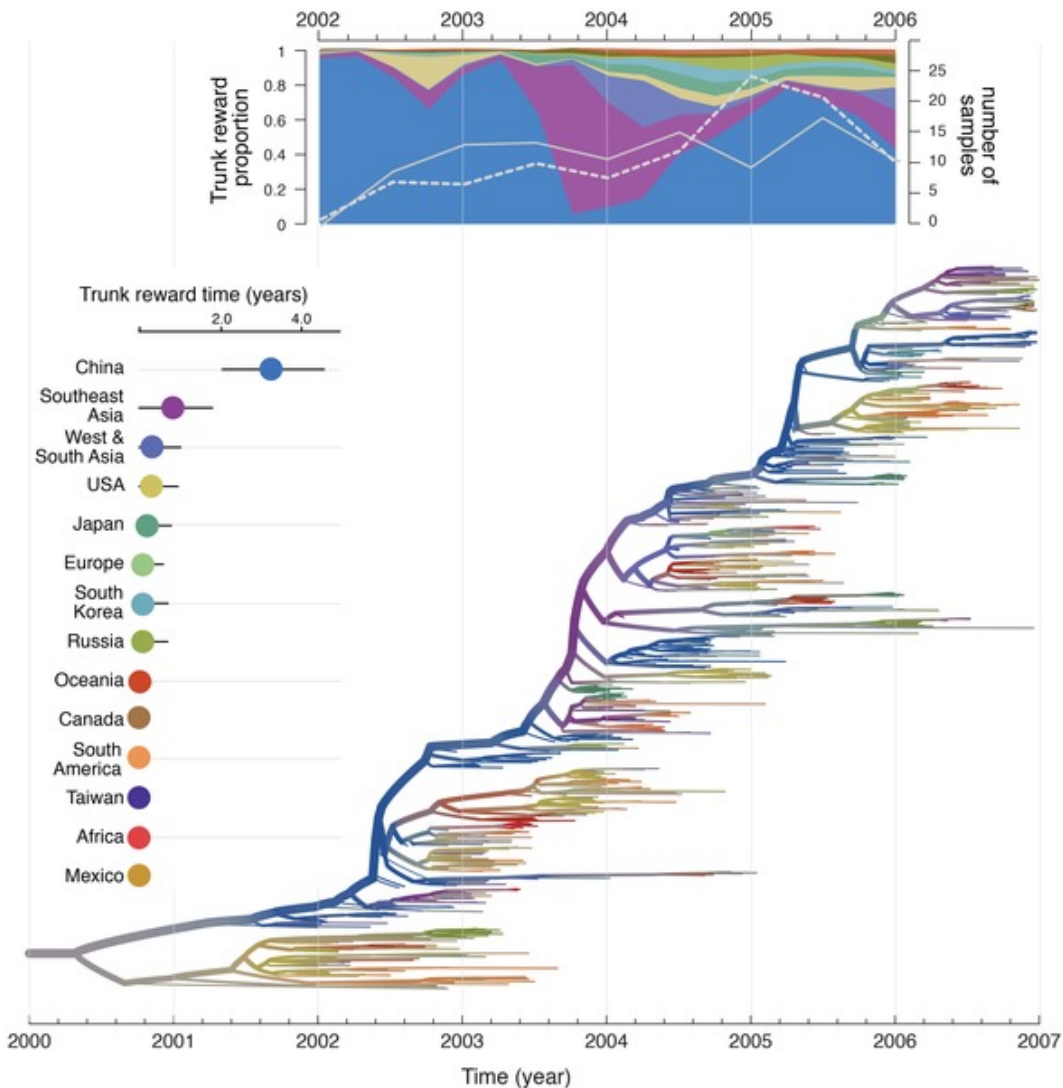
Discrete trait analyses treat geography the same way as substitutions are treated

- Some number of states instead of the 4 nucleotides are considered, but only one position. The number of states corresponds to discrete locations.
- Some transition model between states (can be asymmetric) which are the migration rates.
- Is not a tree prior, as such, a separate tree prior is needed.

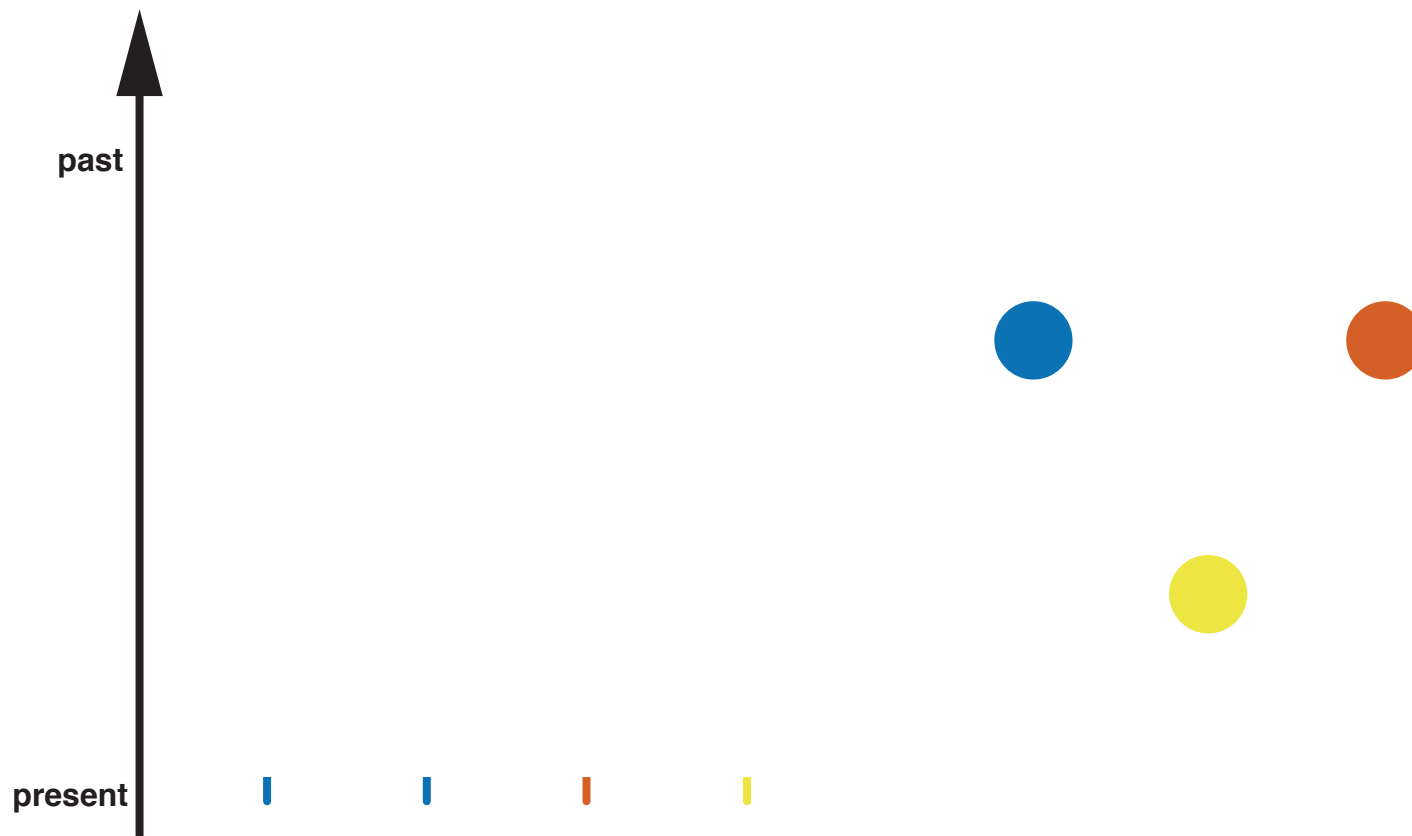
Migration rates can be informed using outside data in a GLM approach

$$\beta \text{ exp} \left[\begin{array}{c} \text{predictor 1} \\ \begin{array}{|c|} \hline a \\ \hline b \\ \hline c \\ \hline d \\ \hline e \\ \hline \end{array} \end{array} \right] + \beta^2 \sigma^2 \left[\begin{array}{c} \text{predictor 2} \\ \begin{array}{|c|} \hline a \\ \hline b \\ \hline c \\ \hline d \\ \hline e \\ \hline \end{array} \end{array} \right] + \beta^3 \sigma^3 \left[\begin{array}{c} \text{predictor 3} \\ \begin{array}{|c|} \hline a \\ \hline b \\ \hline c \\ \hline d \\ \hline e \\ \hline \end{array} \end{array} \right] = \left[\begin{array}{c} \text{Ne or m} \\ \begin{array}{|c|} \hline a \\ \hline b \\ \hline c \\ \hline d \\ \hline e \\ \hline \end{array} \end{array} \right]$$

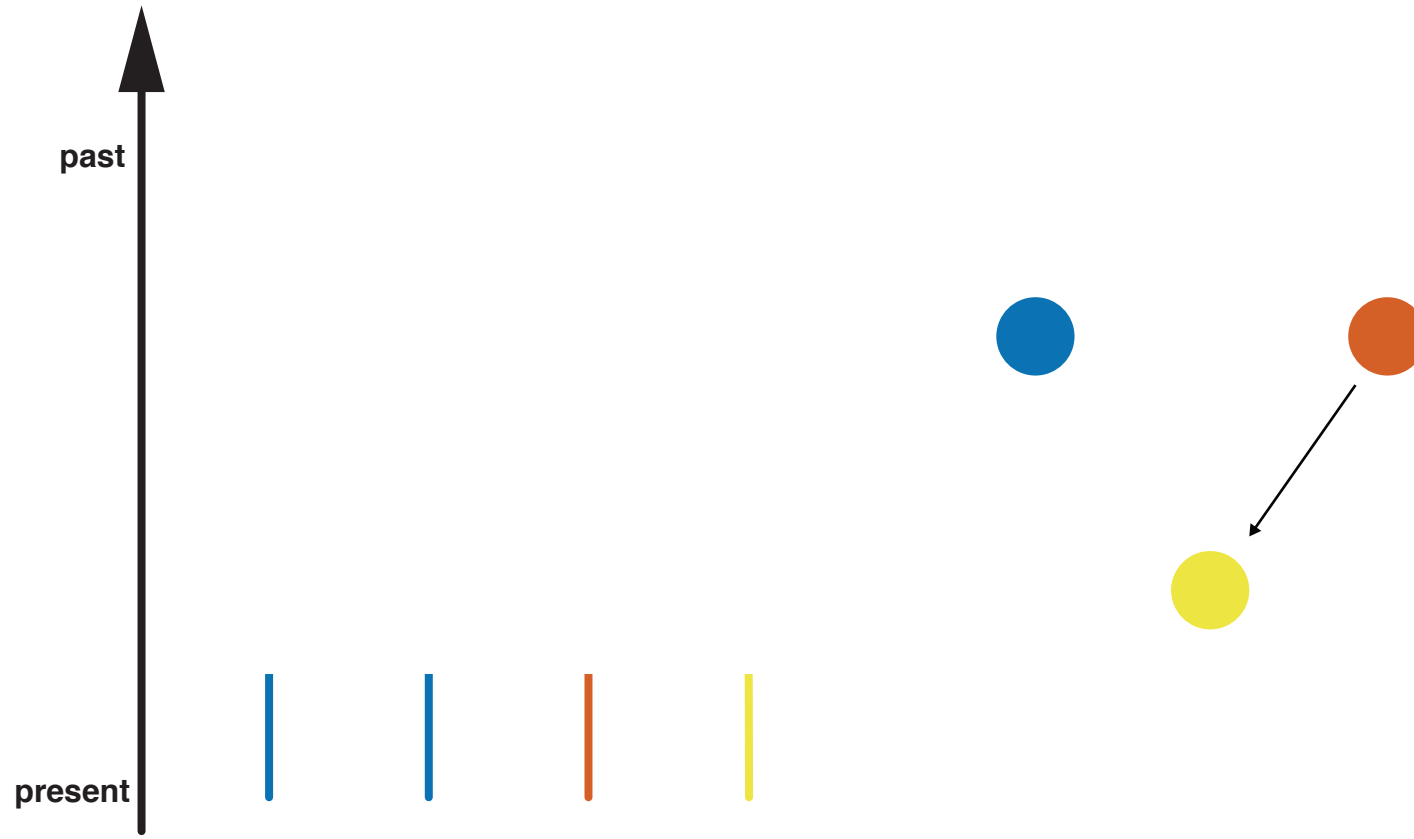
Migration rates in the DTA model can be informed by transportation Data



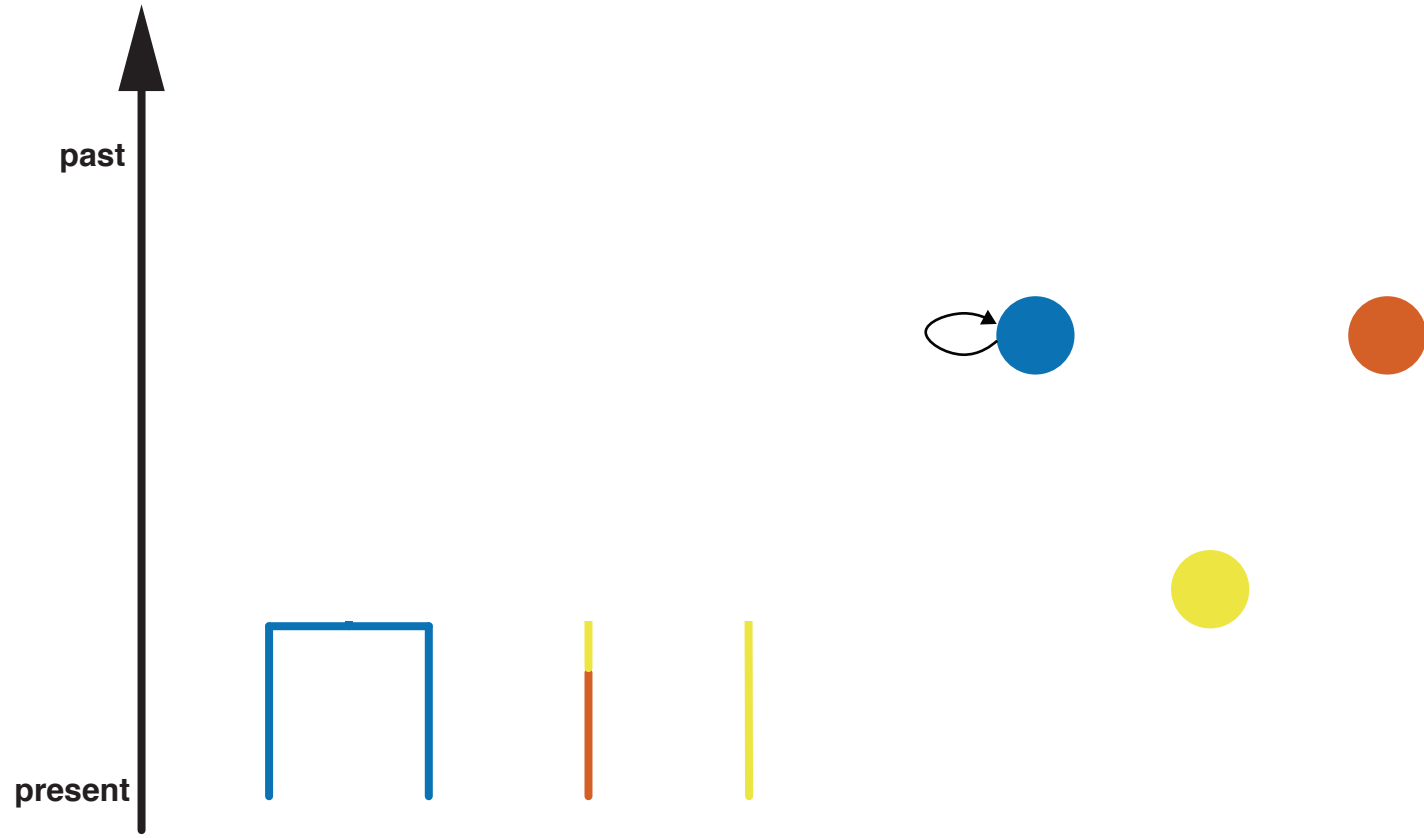
We track the ancestry of 4 viruses isolated from 4 people
in 3 locations from present to past



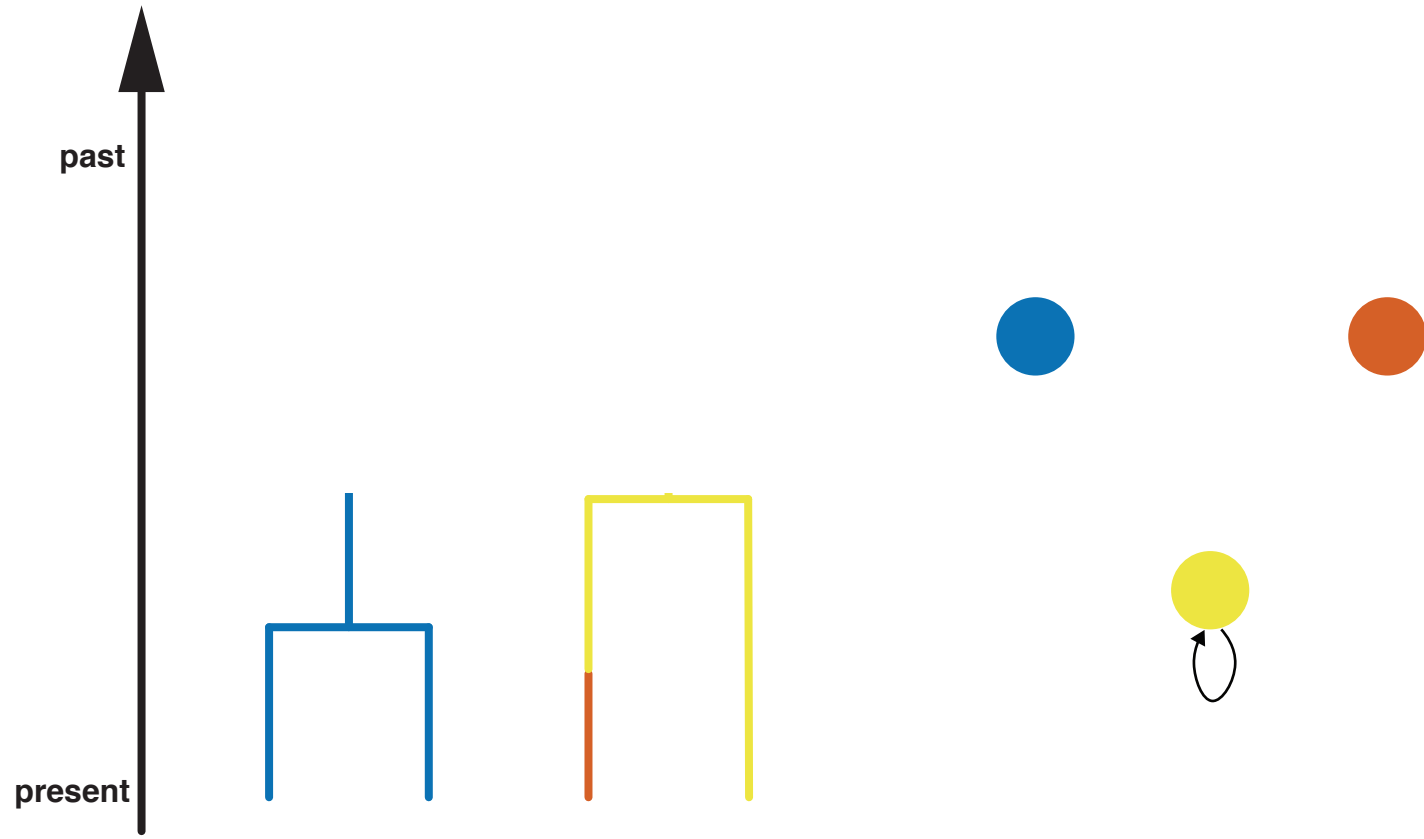
Lineages can have originated from a different location at a rate given by the migration rate.



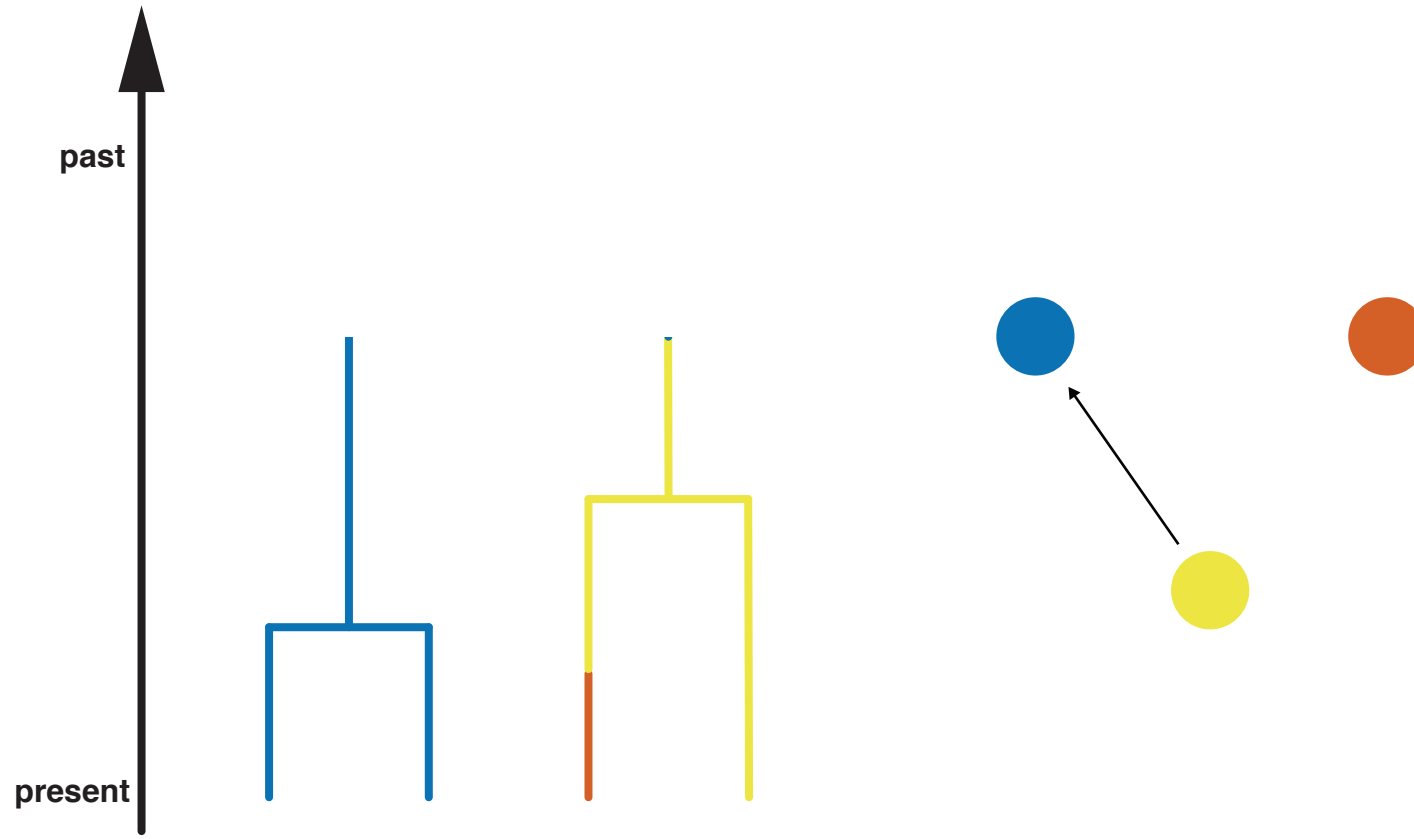
Two lineages can share a common ancestor at a rate
inverse proportional to the N_e of location blue



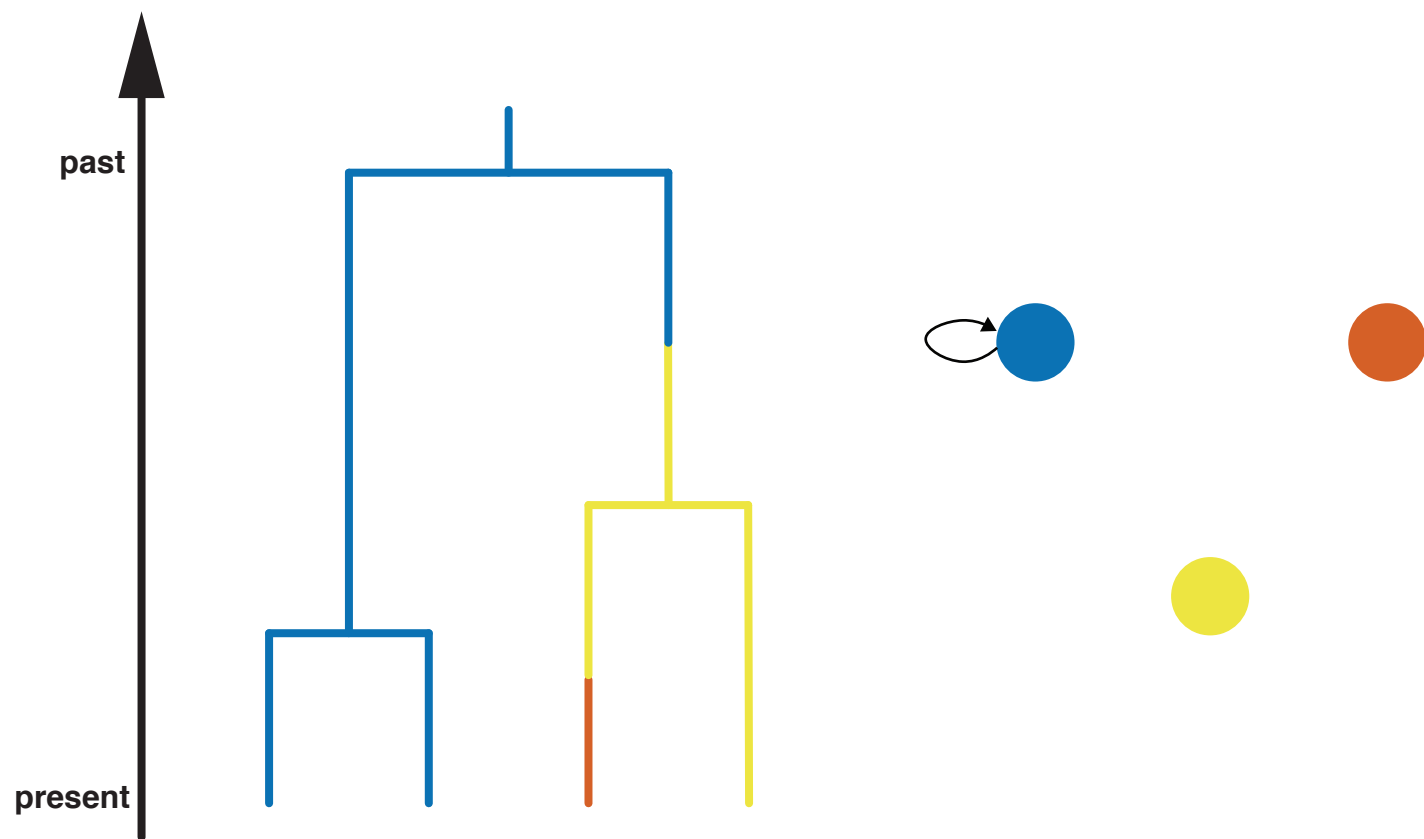
Two lineages can share a common ancestor at a rate
inverse proportional to the N_e of location yellow



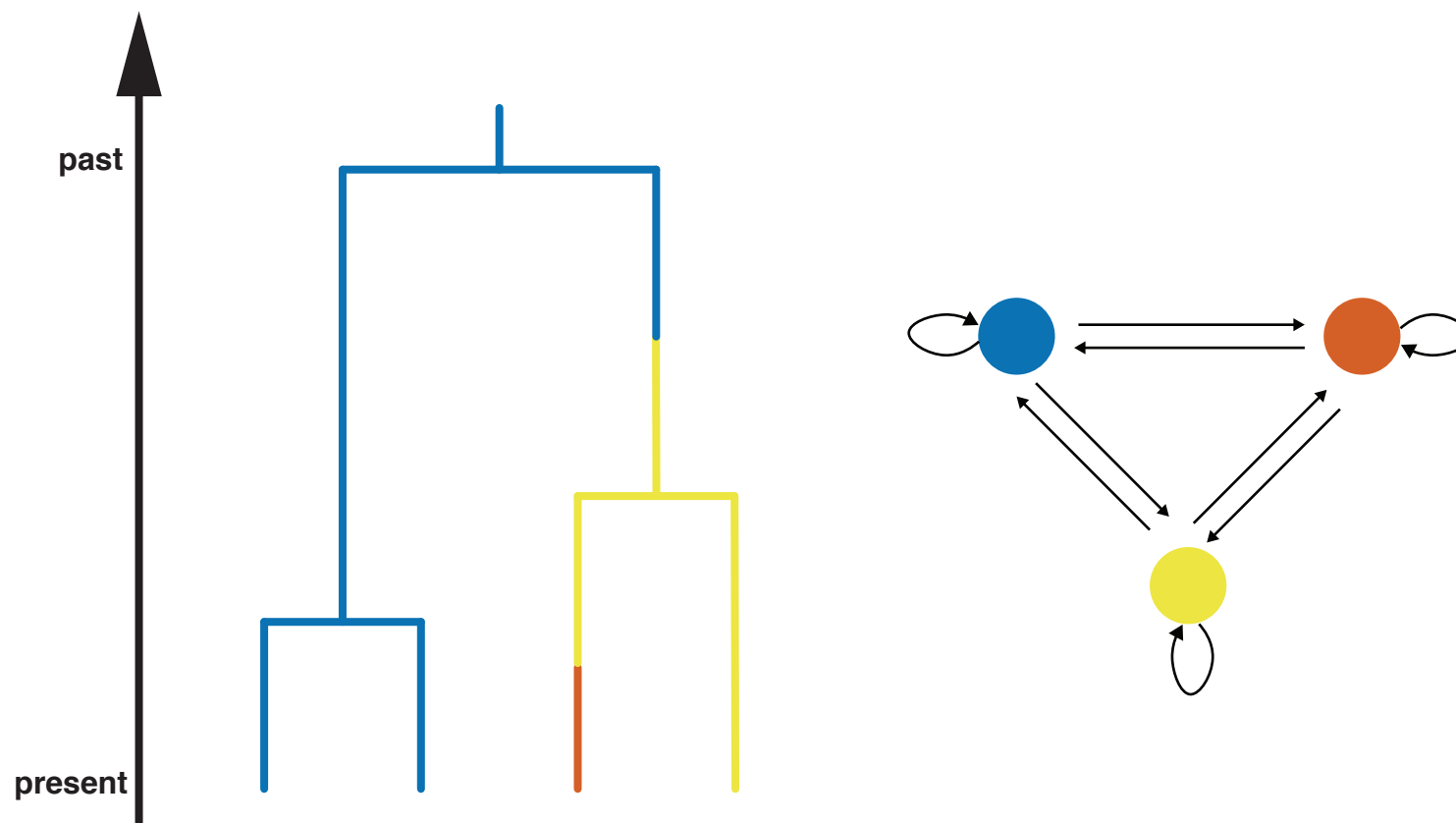
Lineages can have originated from a different location at a rate given by the migration rate.



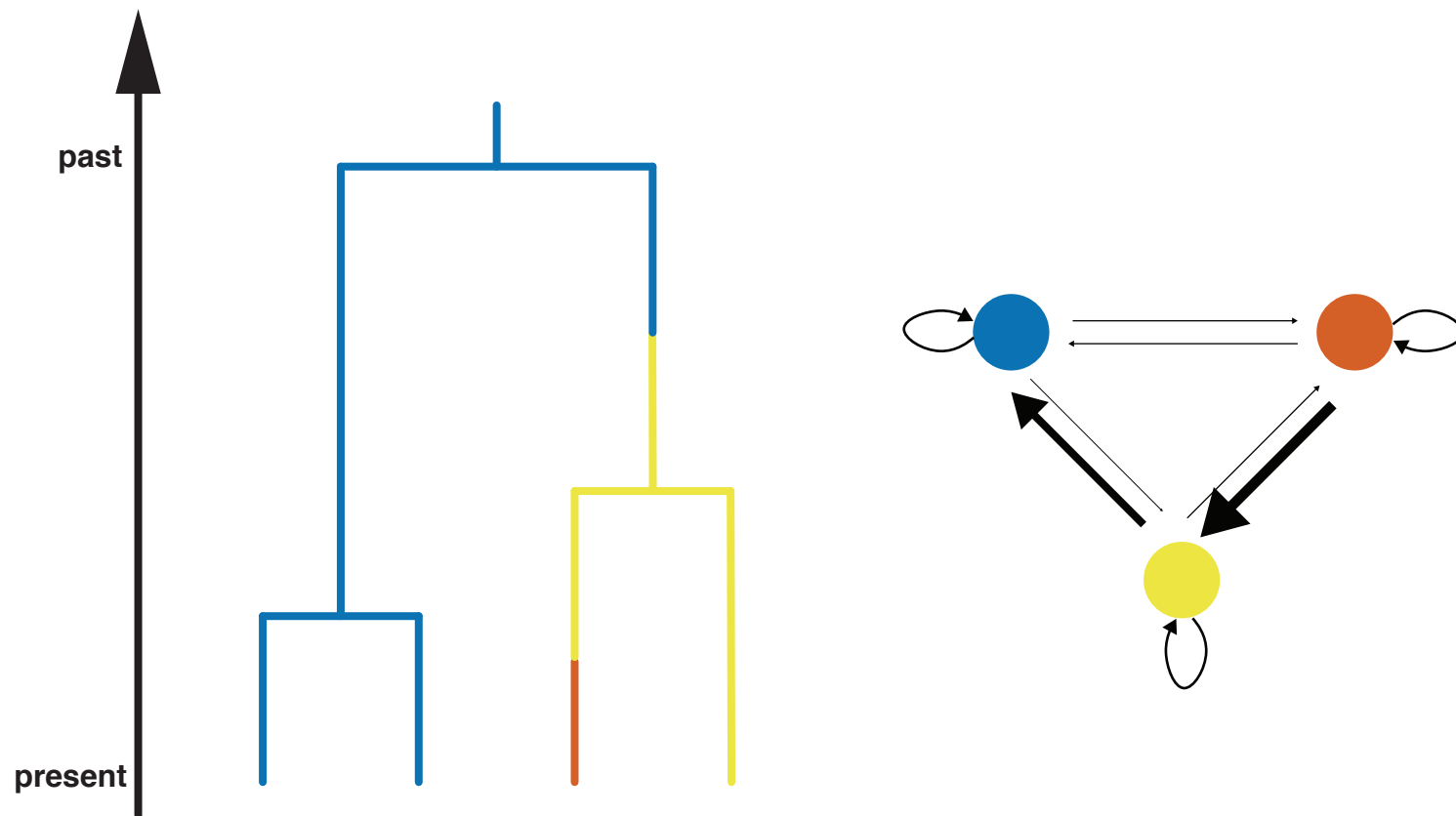
And then coalesce again in location blue, reaching the root of the tree



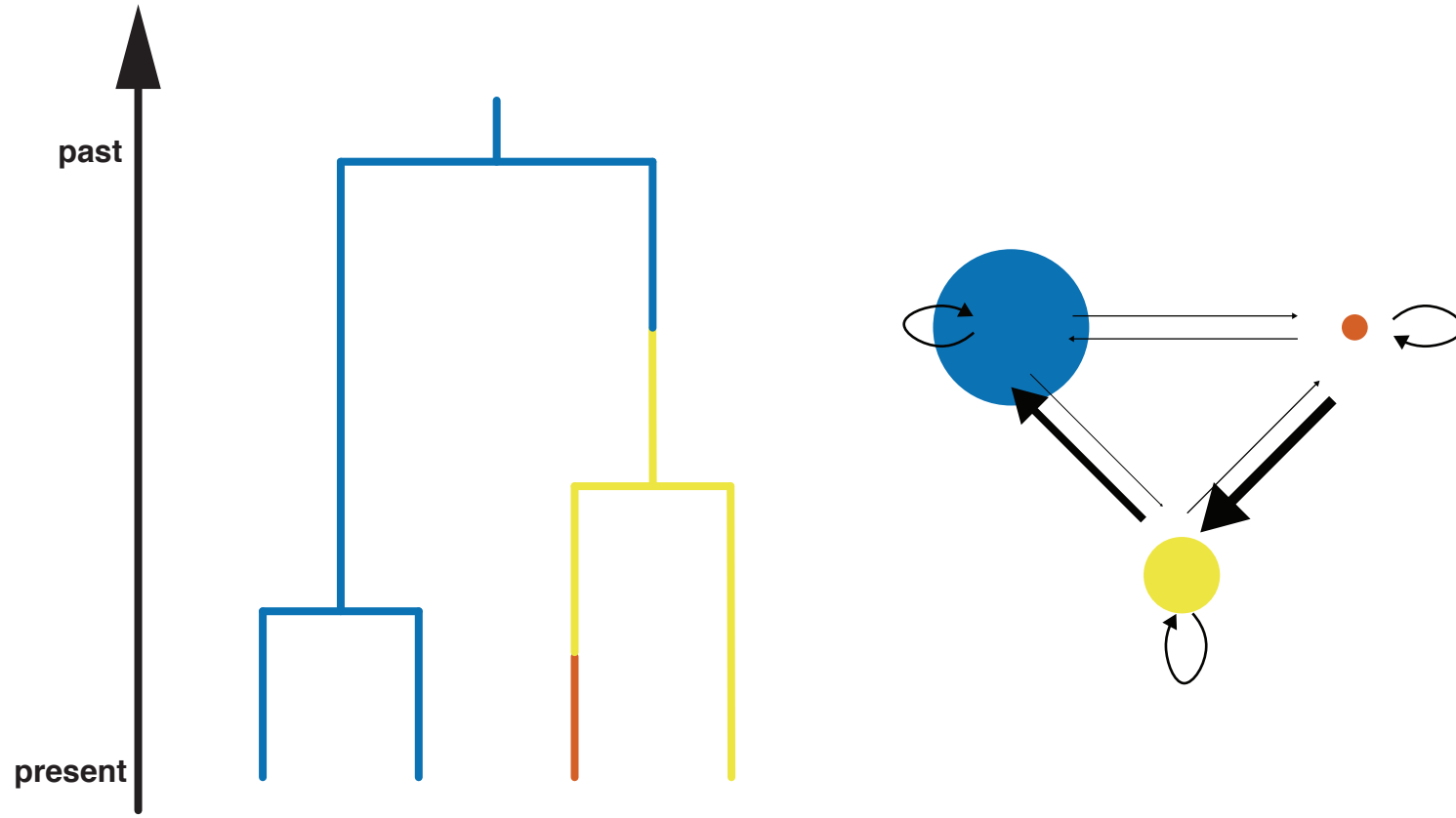
The structured coalescent models how lineages coalesce within and migrate between locations



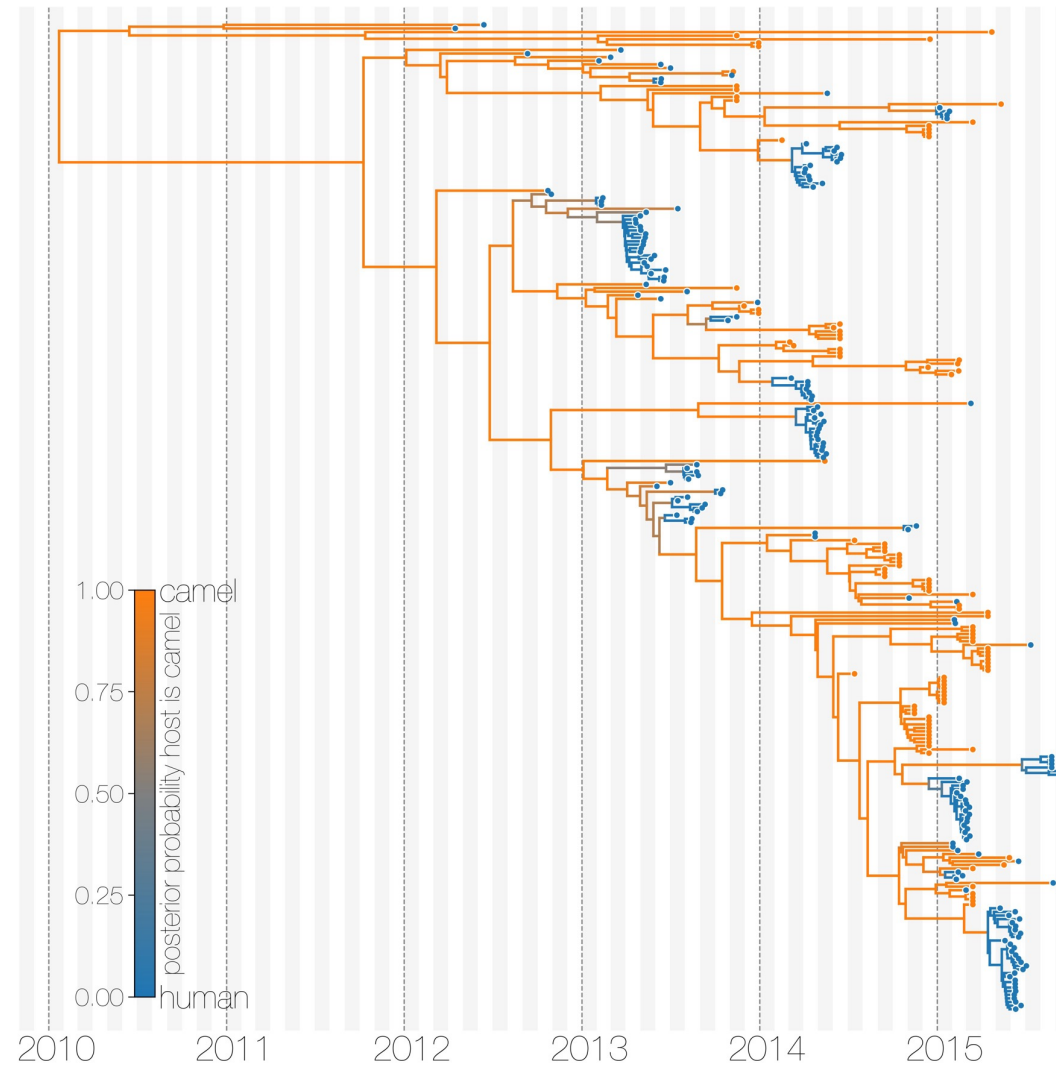
This allows us to estimate how migration rates between locations differ



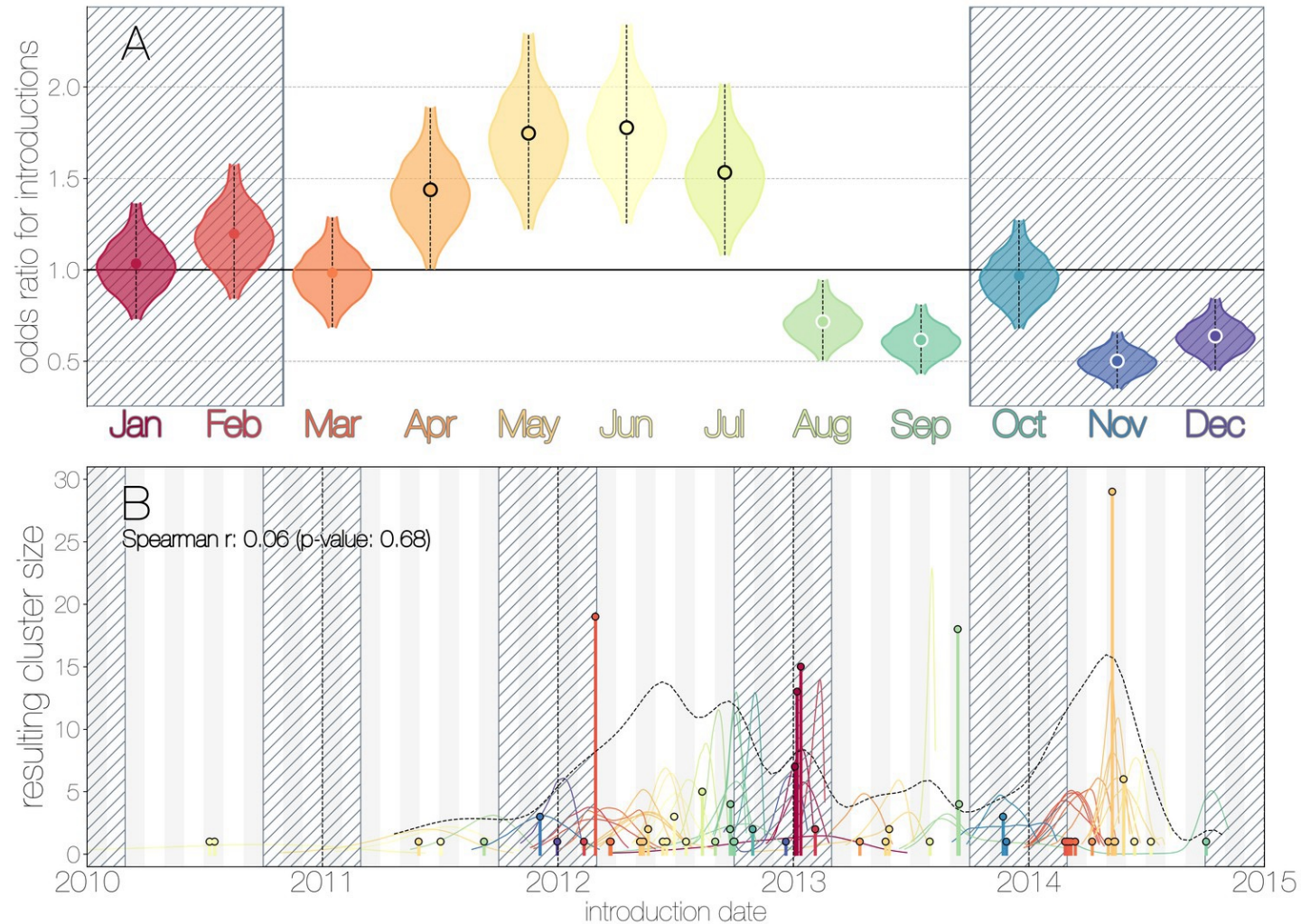
Or how the number of infected individuals differs between location



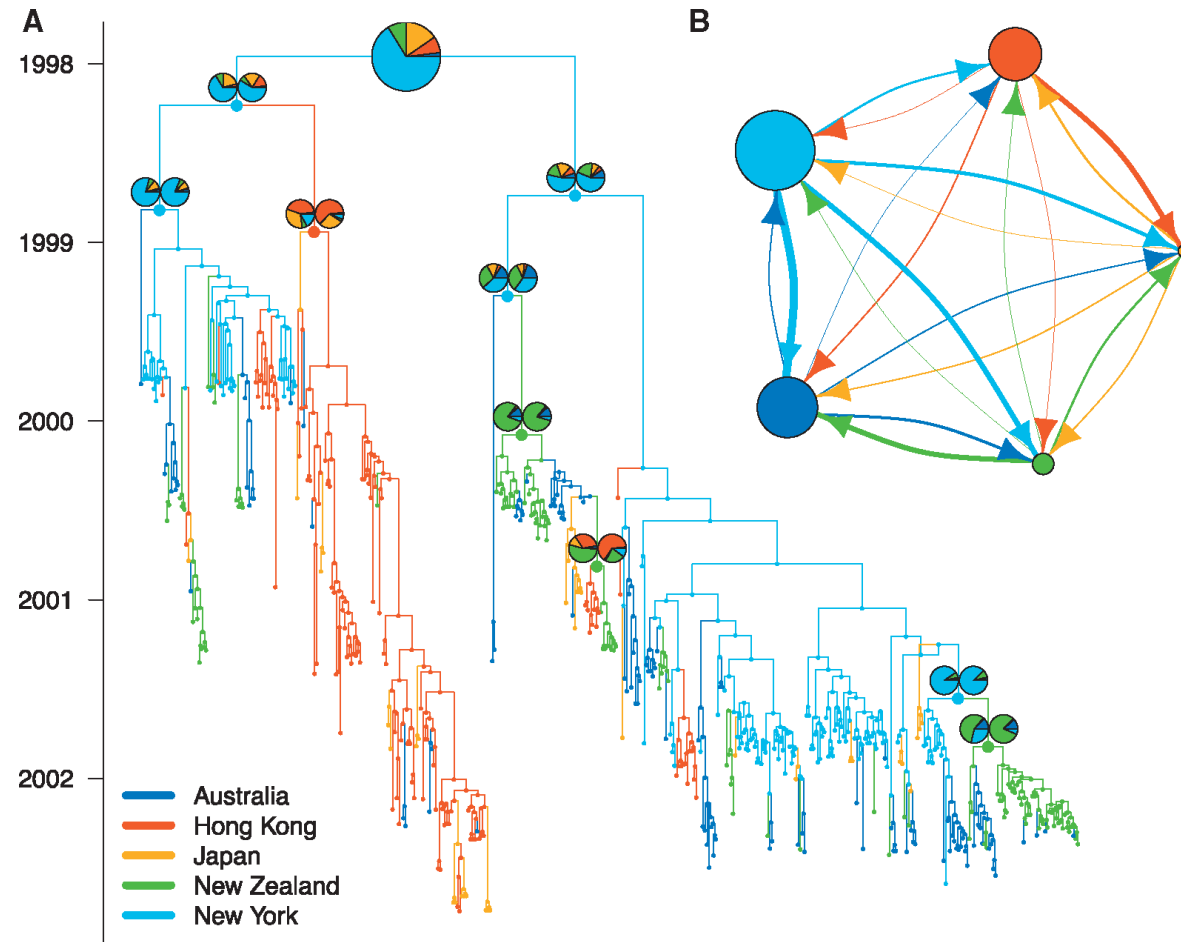
MERS sequences cluster by host of isolation



Seasonal trends in camel to human transmissions



MASCOT approximates the structured coalescent to dramatically improve efficiency



Migration rates and effective population sizes can be informed using outside data in a GLM approach

$$\beta \text{ exp} \left[\begin{array}{c} \text{predictor 1} \\ \begin{array}{|c|} \hline a \\ \hline b \\ \hline c \\ \hline d \\ \hline e \\ \hline \end{array} \end{array} + \beta^2 \sigma^2 \begin{array}{c} \text{predictor 2} \\ \begin{array}{|c|} \hline a \\ \hline b \\ \hline c \\ \hline d \\ \hline e \\ \hline \end{array} \end{array} + \beta^3 \sigma^3 \begin{array}{c} \text{predictor 3} \\ \begin{array}{|c|} \hline a \\ \hline b \\ \hline c \\ \hline d \\ \hline e \\ \hline \end{array} \end{array} \right] = \begin{array}{c} \text{Ne or m} \\ \begin{array}{|c|} \hline a \\ \hline b \\ \hline c \\ \hline d \\ \hline e \\ \hline \end{array} \end{array}$$

In 2016 and 2017, there was an Influenza A/H5N8 outbreak in France

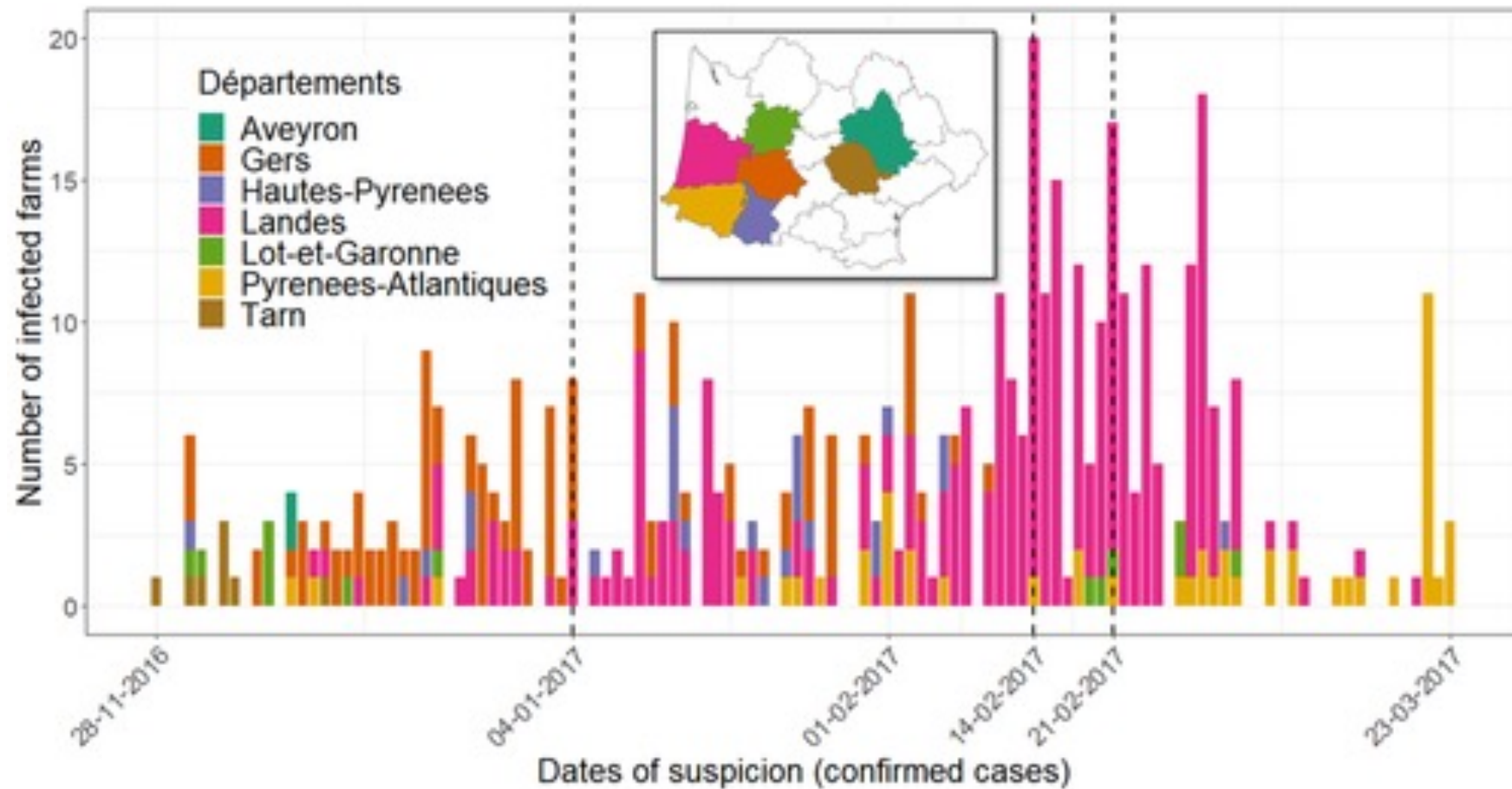
French Foie Gras Makers Fear the Worst as Bird Flu Toll Rises

The bird flu epidemic in southwest France, home to most duck and foie gras producers, has led to more than three million poultry being killed.

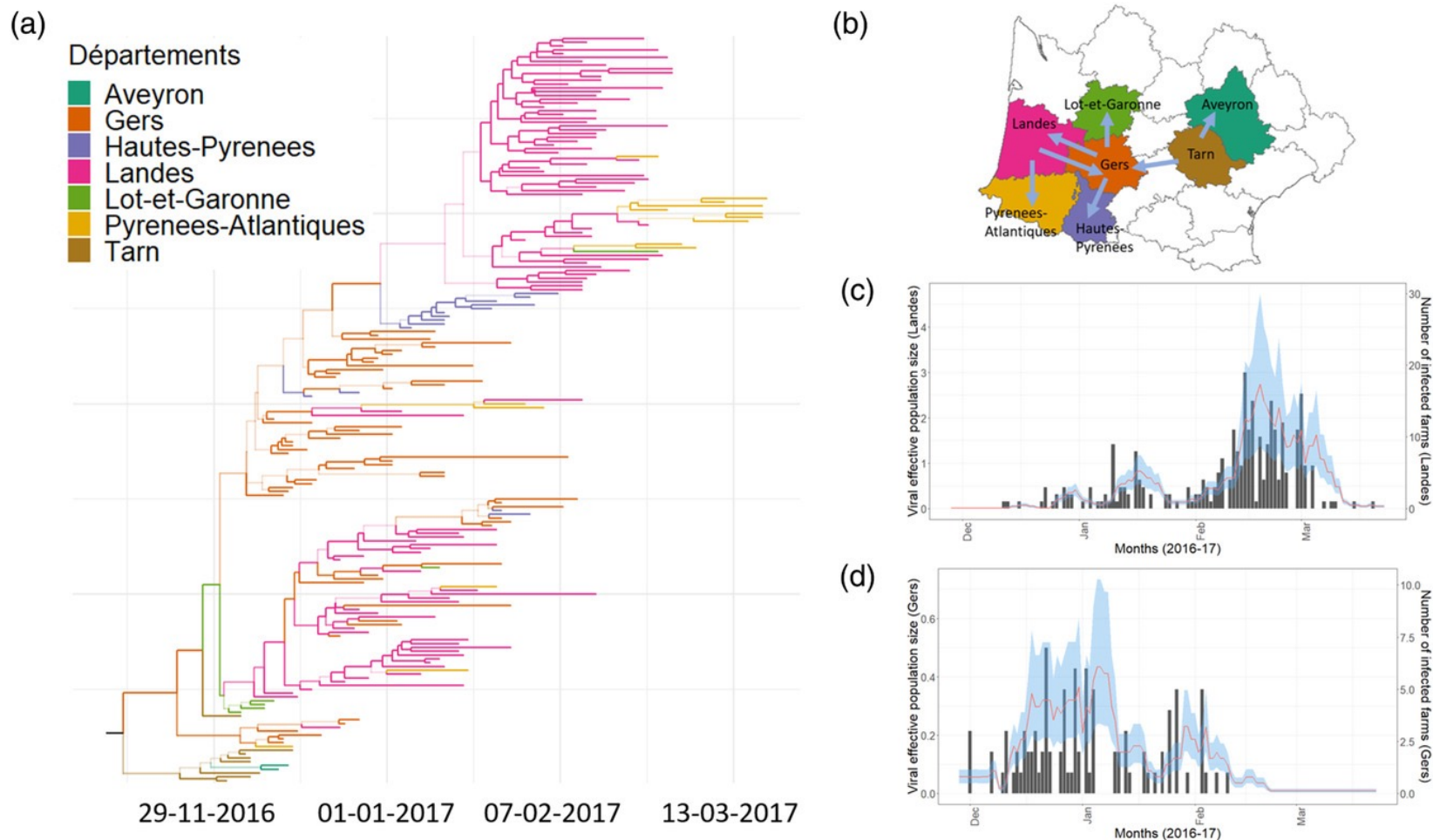


<https://www.nbcnews.com/news/world/french-foie-gras-makers-fear-worst-bird-flu-toll-rises-n721276>

Cases first appeared in the department (state) of tarn and H5N8 was soon detected in other regions as well

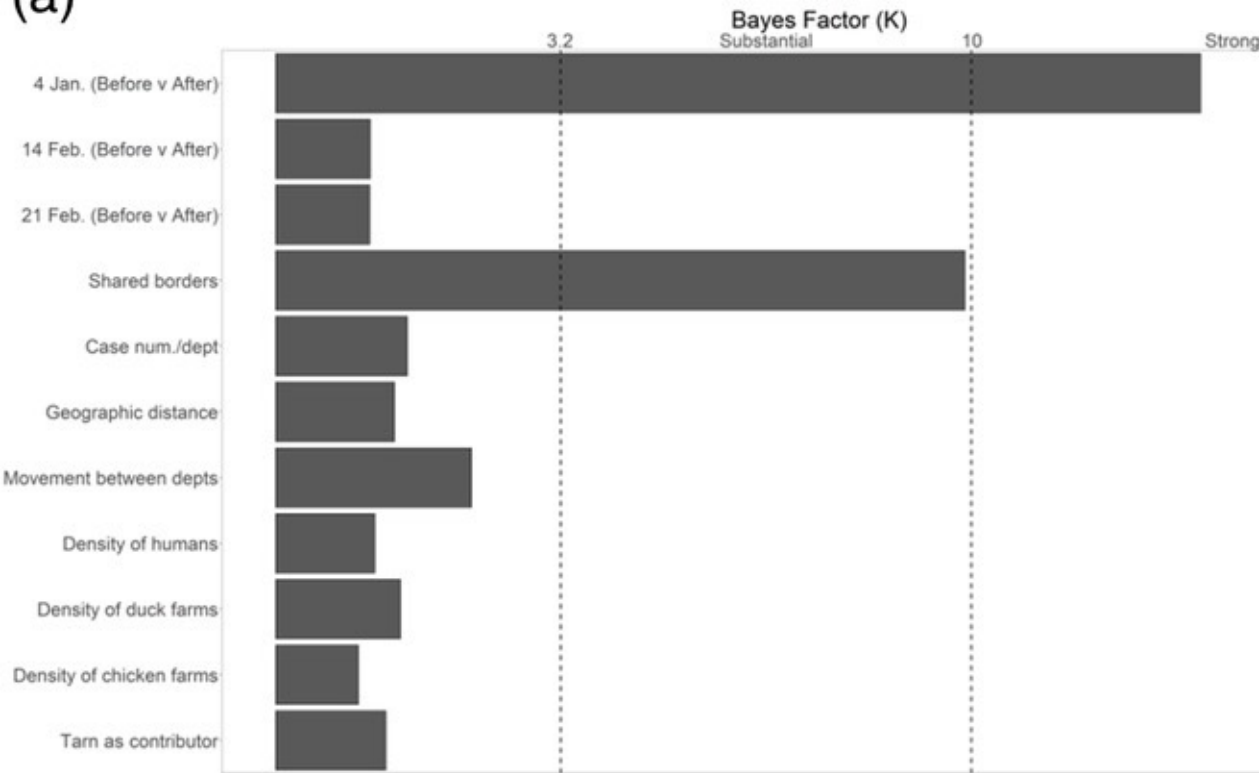


MASCOT-GLM reconstructs Tarn to be the most likely source, consistent with case data

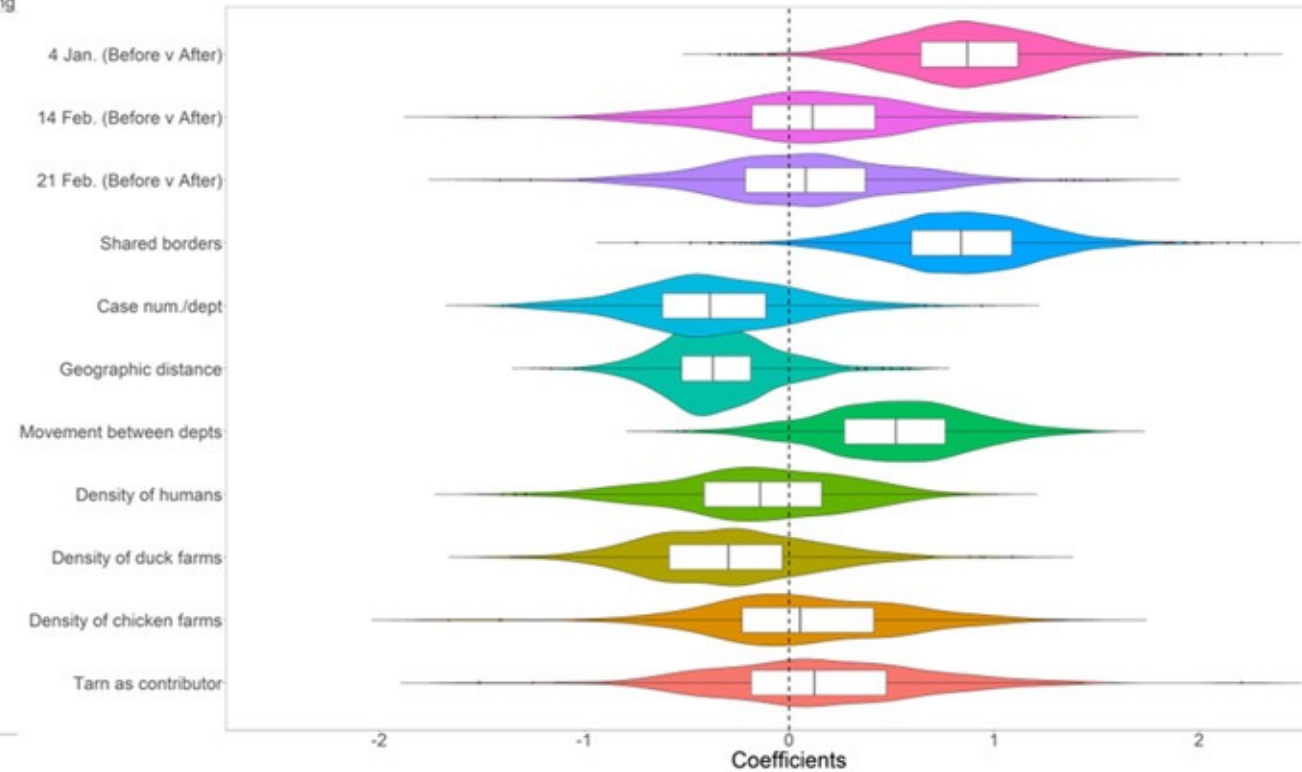


Evidence that the initial culling of birds reduced migration between location

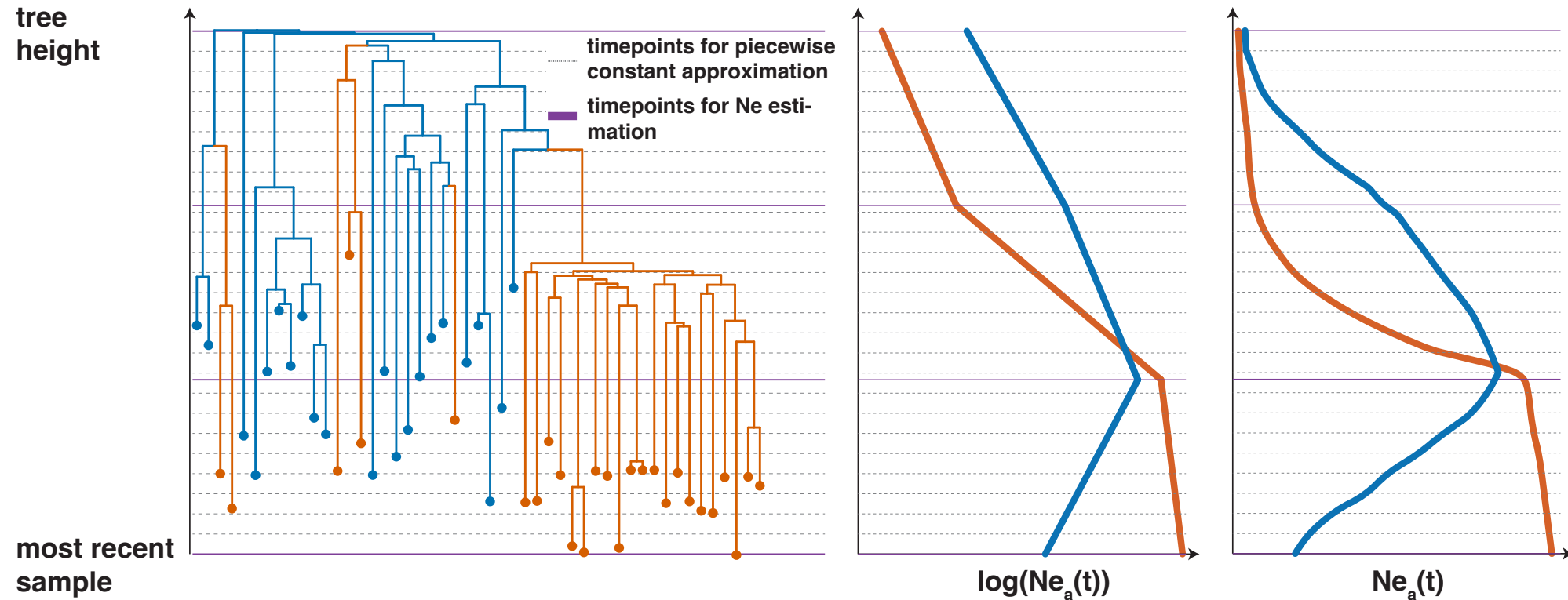
(a)



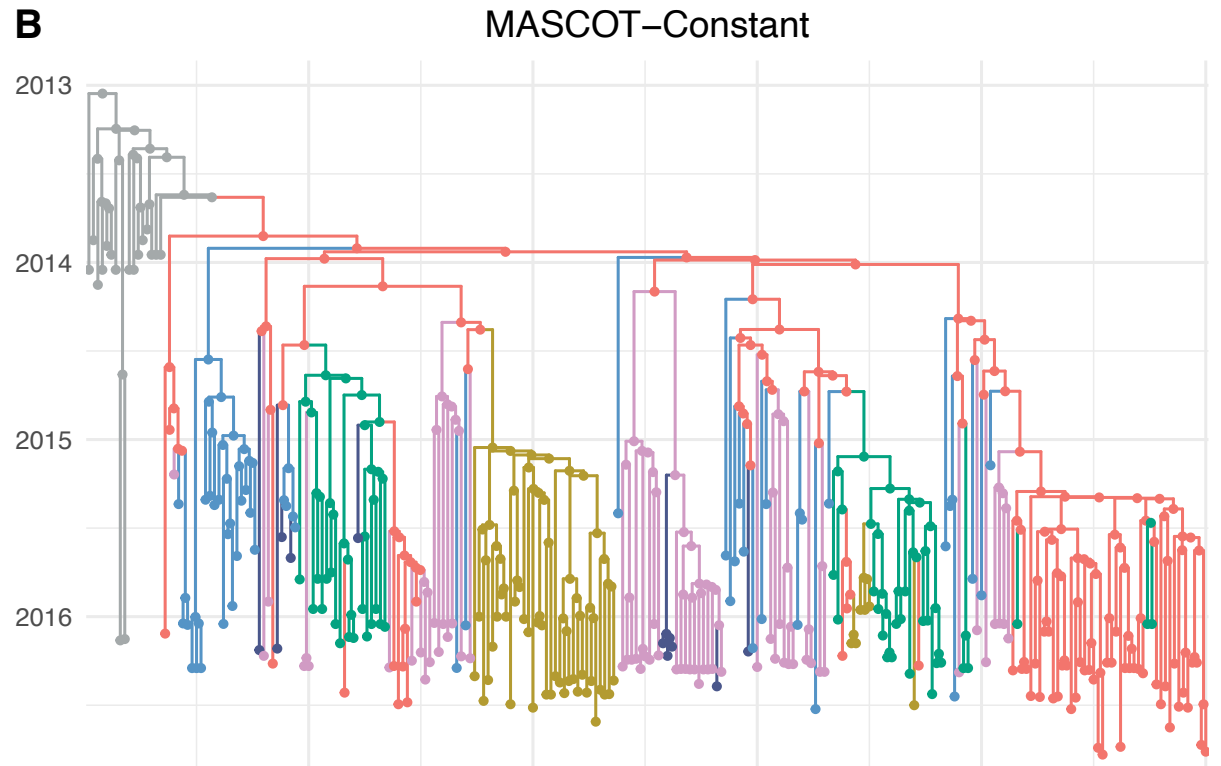
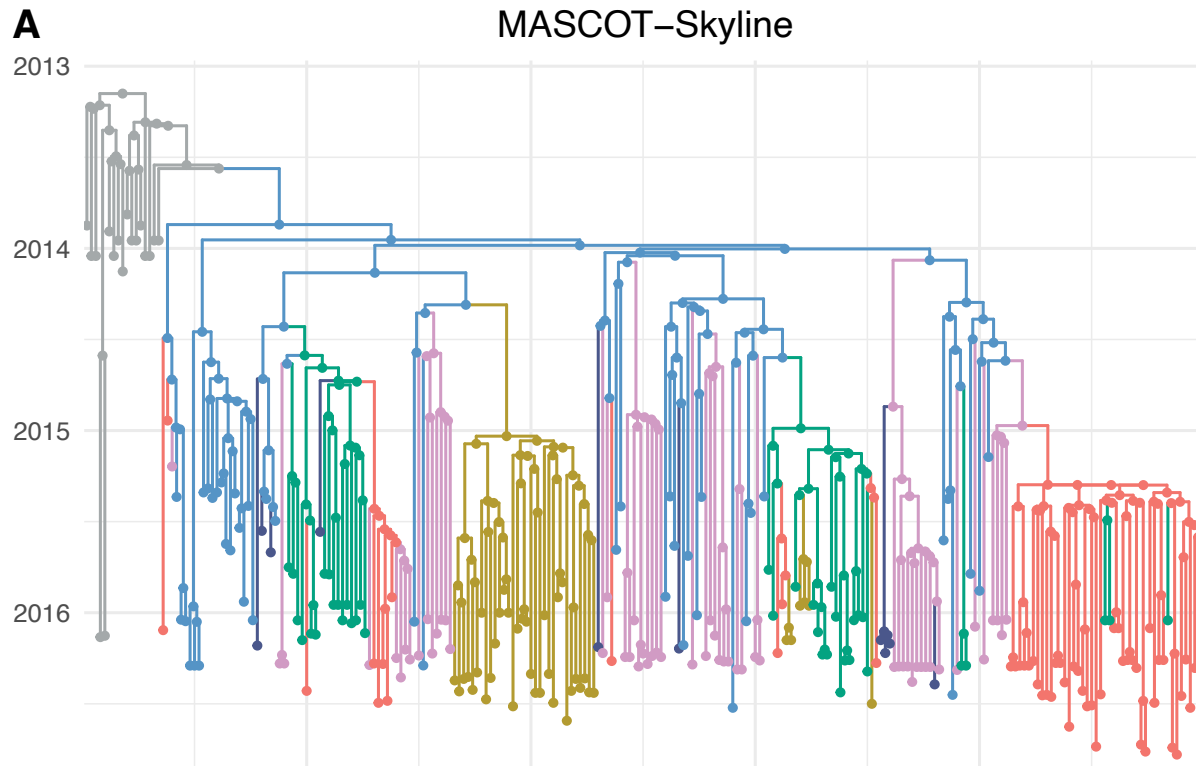
(b)



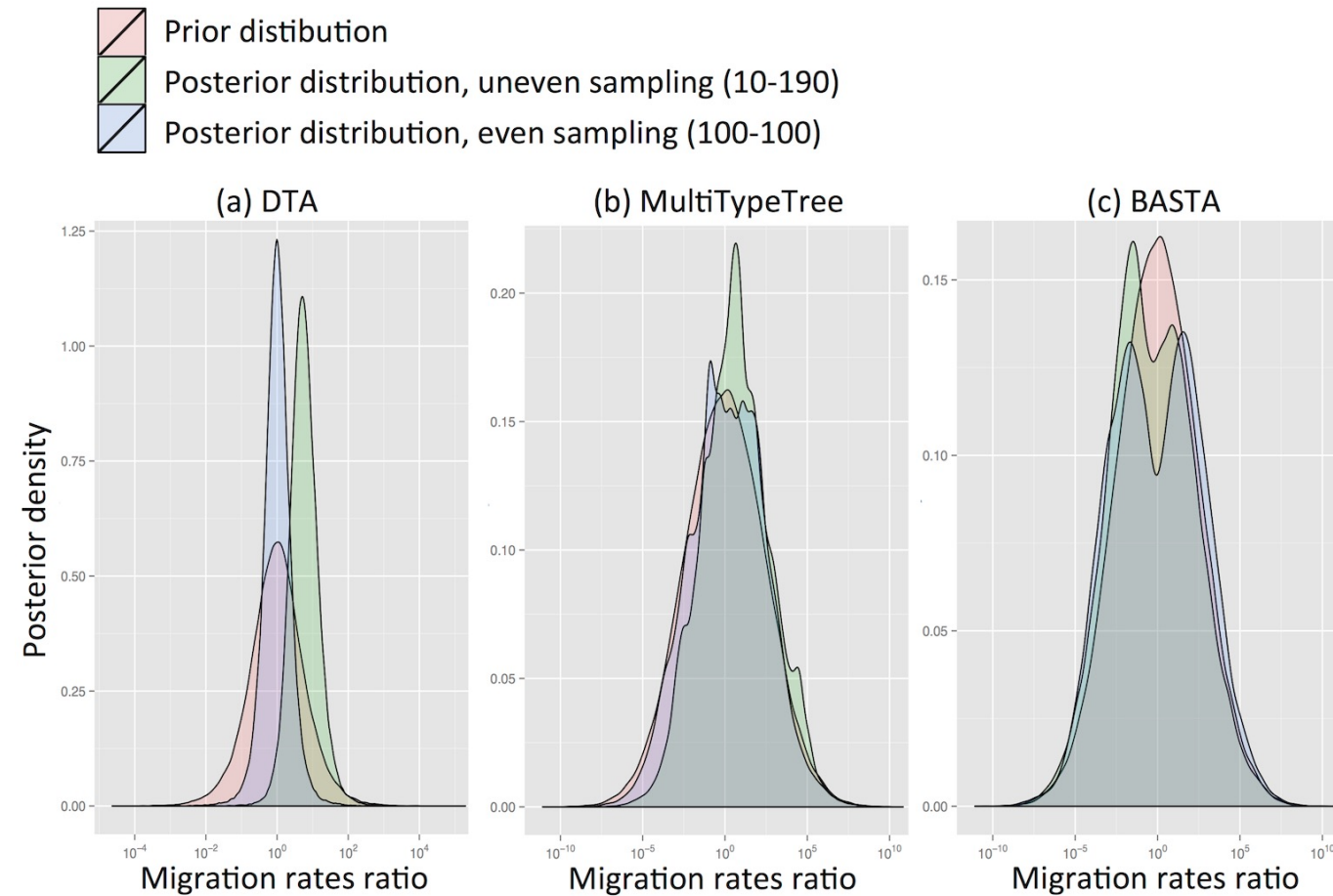
Effective population size dynamics can be modelled using MASCOT-Skyline



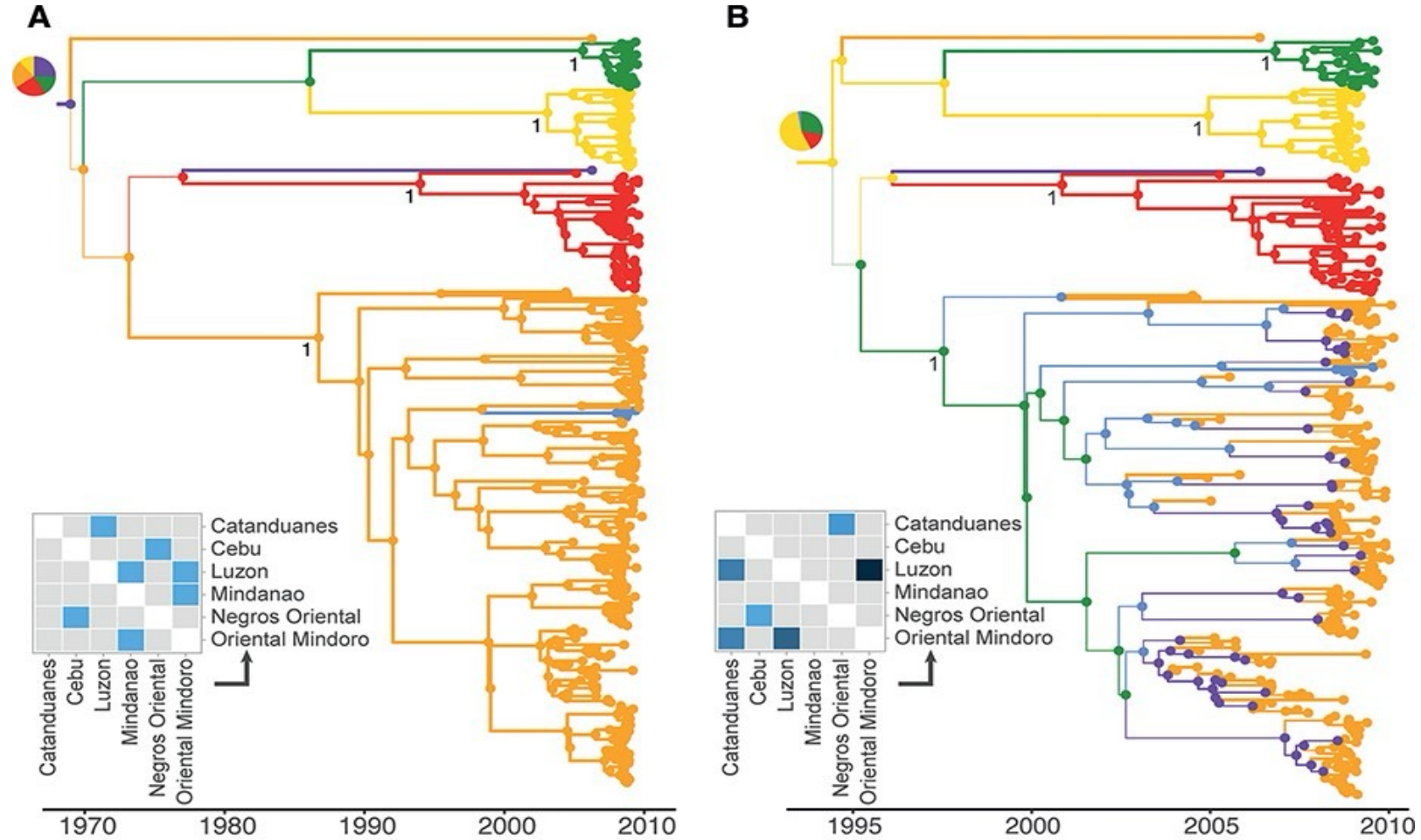
Accounting for population dynamics in structured coalescent can impact ancestral state reconstruction



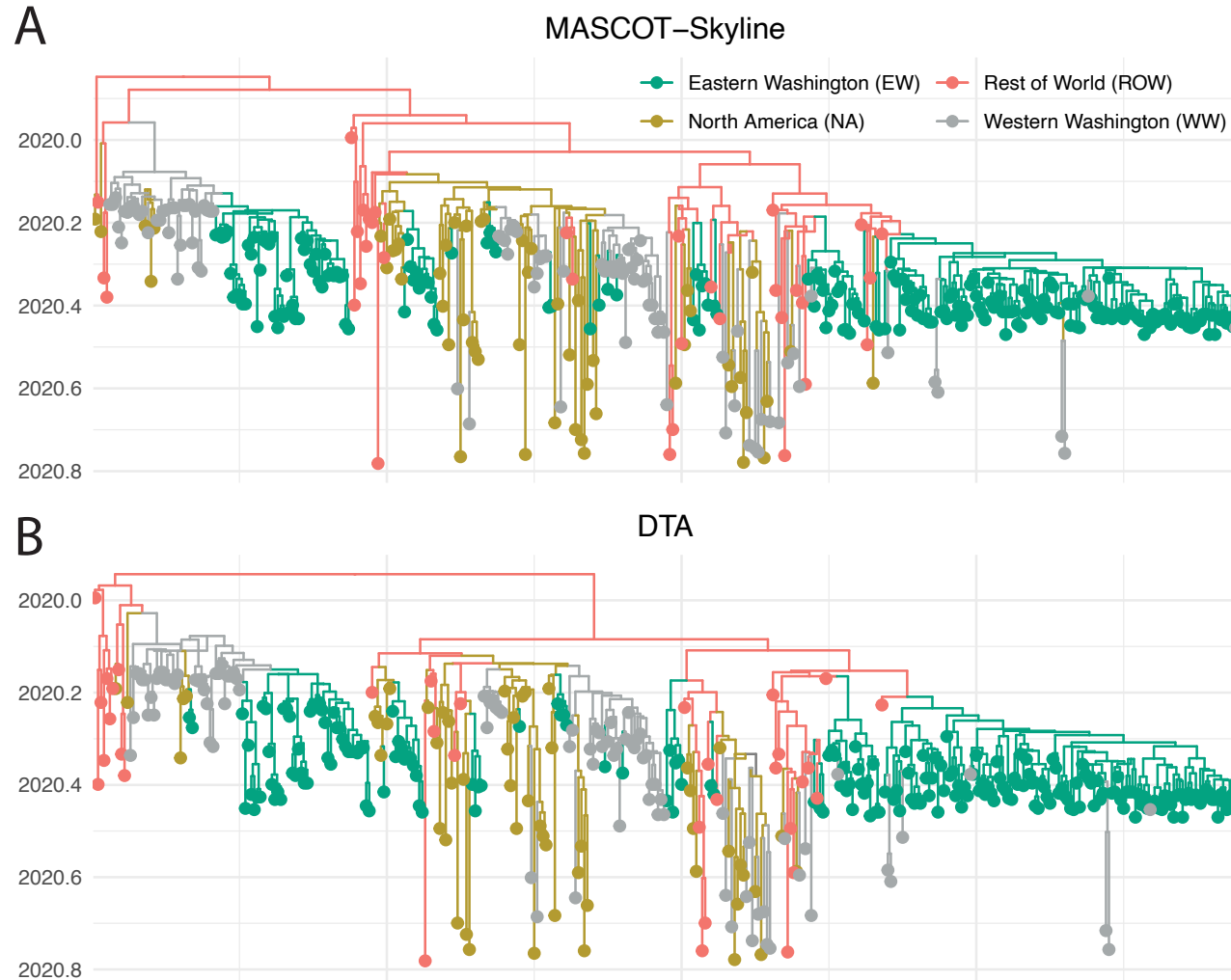
Biases in Structured phylodynamics



DTA typically returns a rather parsimonious ancestral state reconstruction

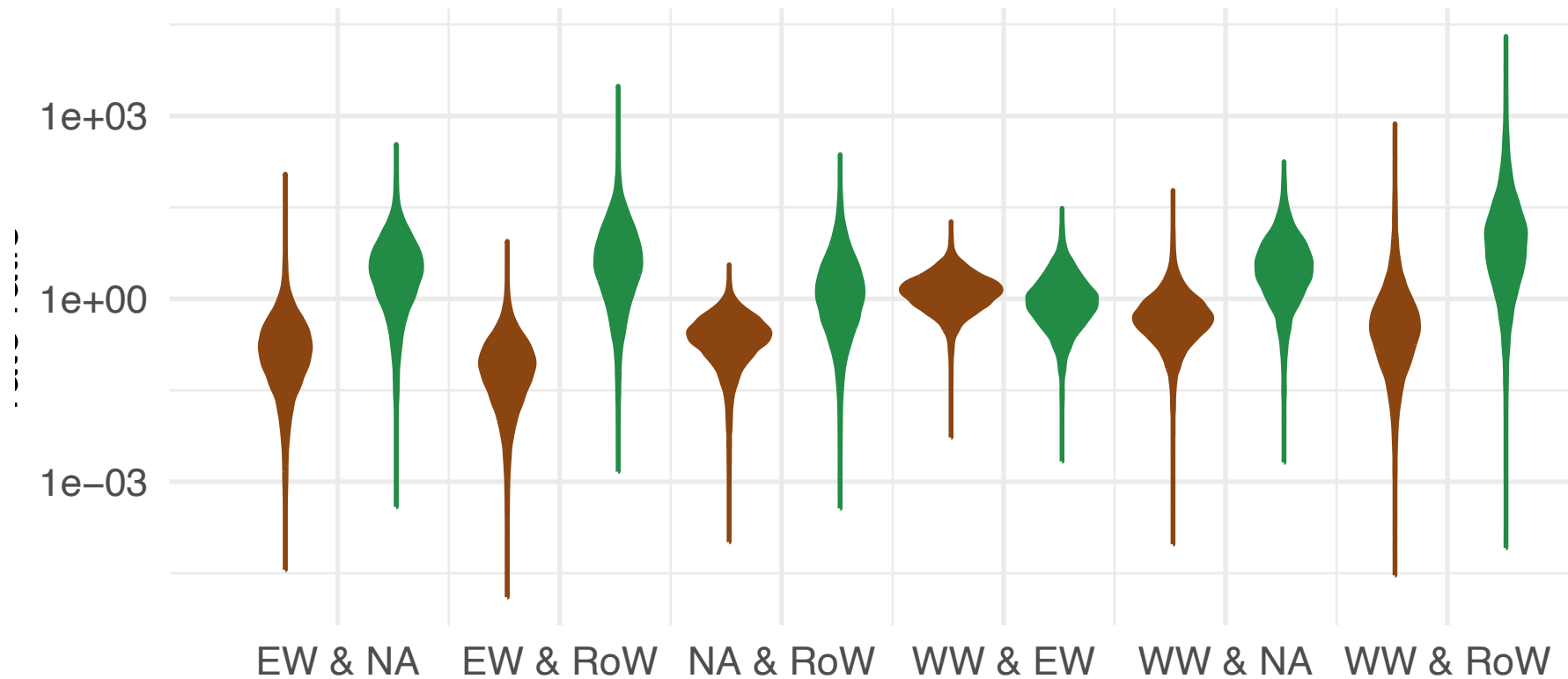


Different methods can interpret very similar observations differently



Different methods can interpret very similar observations differently

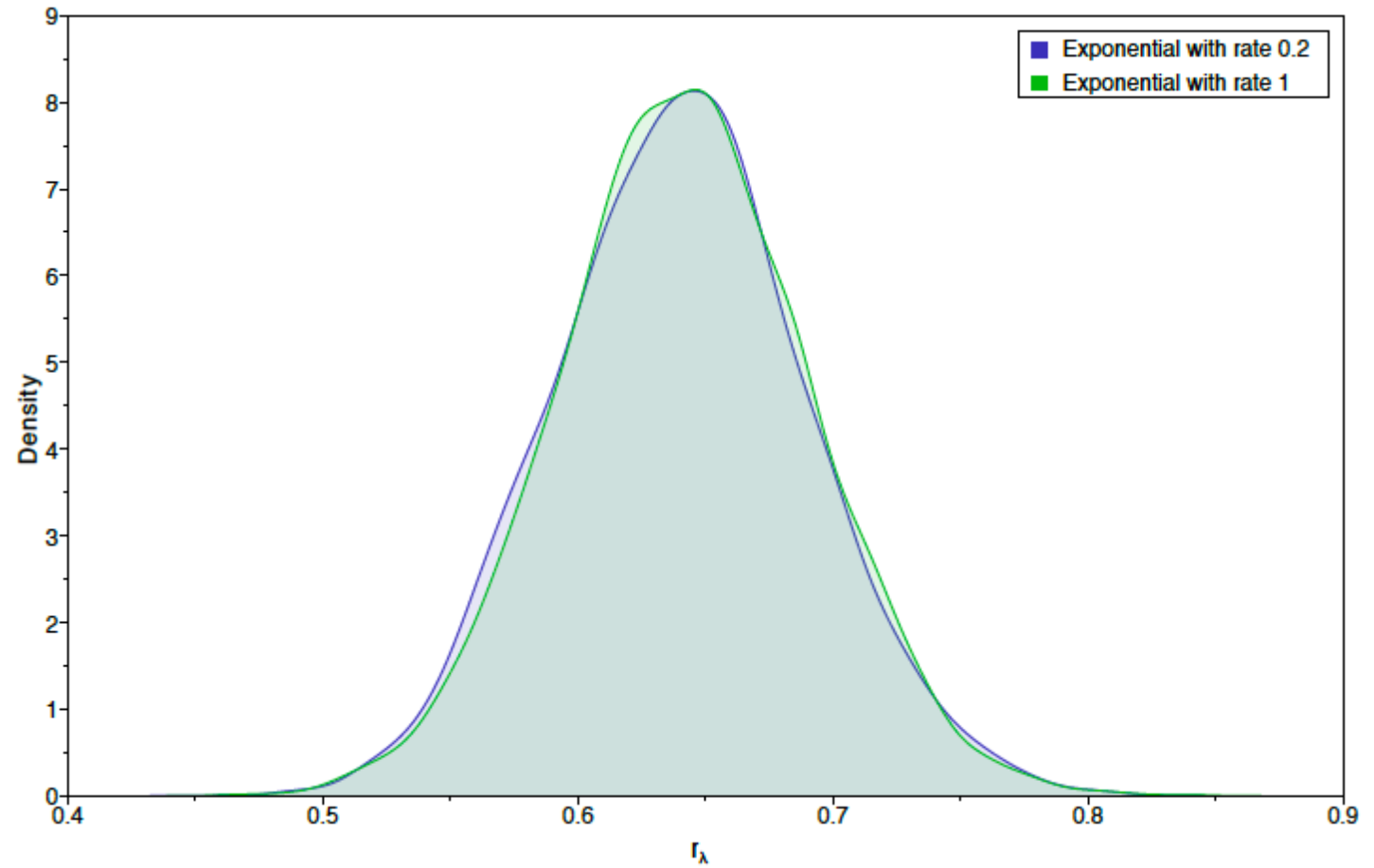
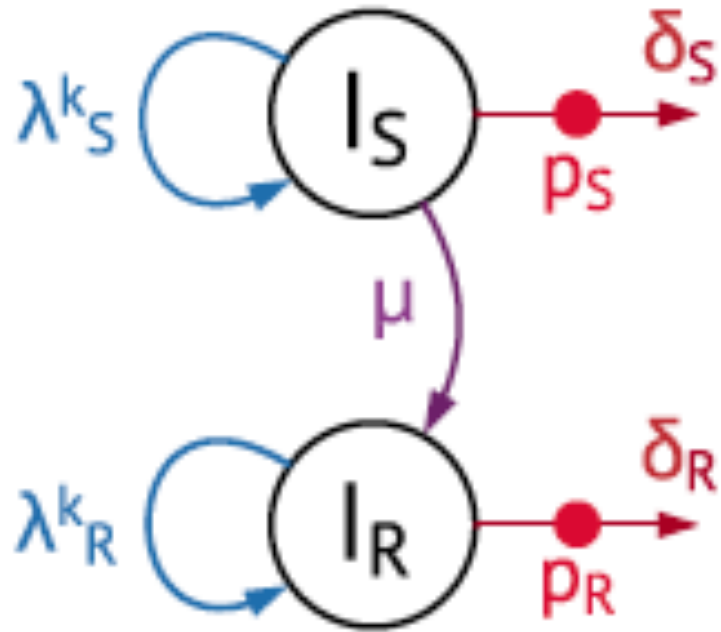
C



Variations

- Typed or untyped trees, i.e. trees including the migration history or not (BDMM)
- Fixed types at the tips (BDMM) or estimated (MSBD)
- Extended multi-type birth-death model with birth rate between demes $\lambda_{i,j} \forall i \neq j$ (BDMM, XML only)

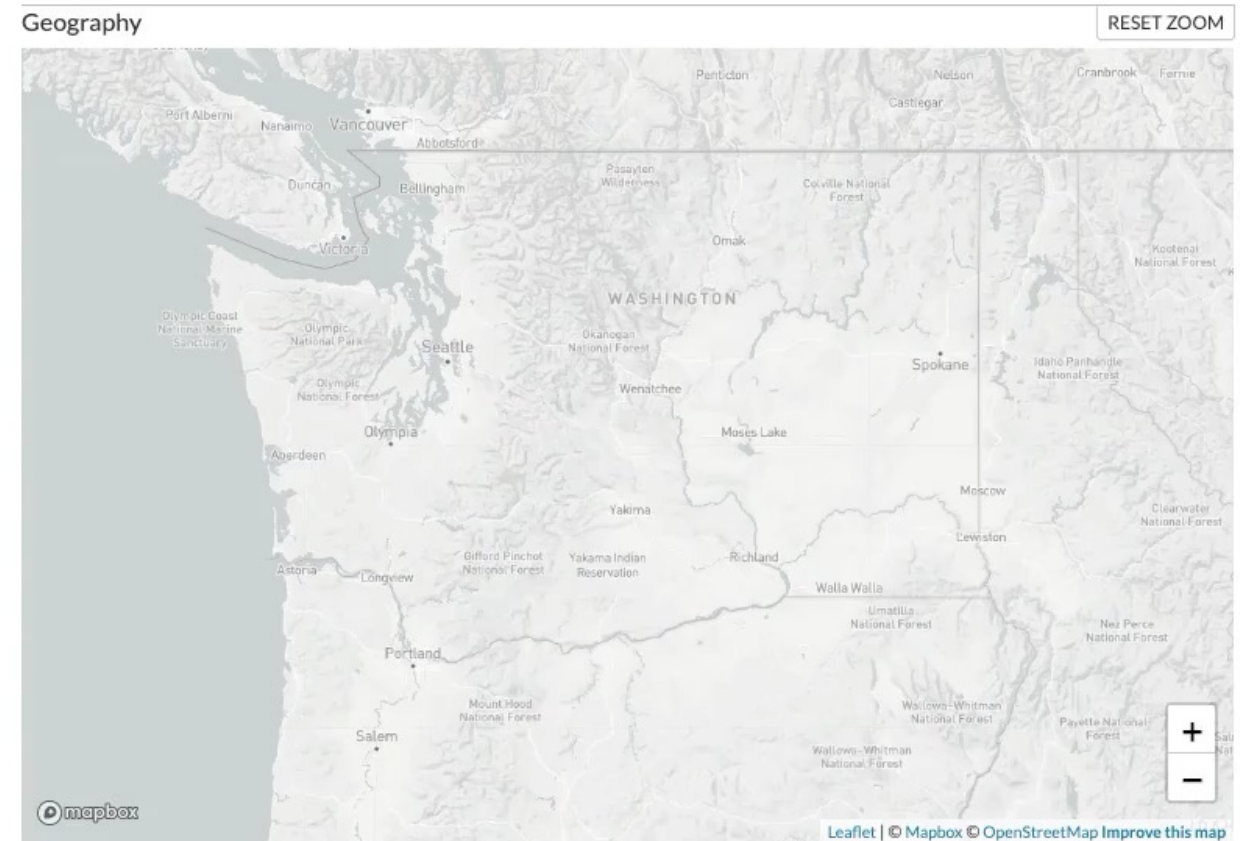
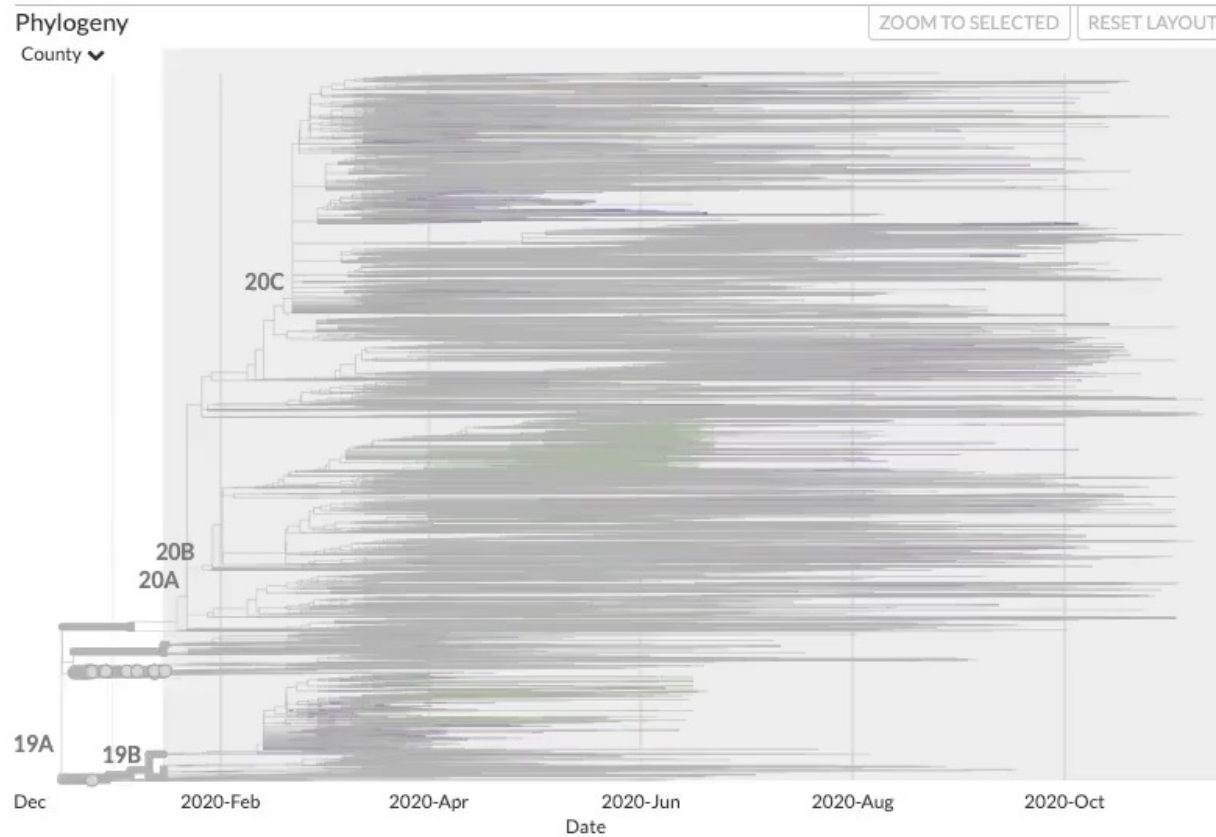
Example: tuberculosis



Pros / Cons

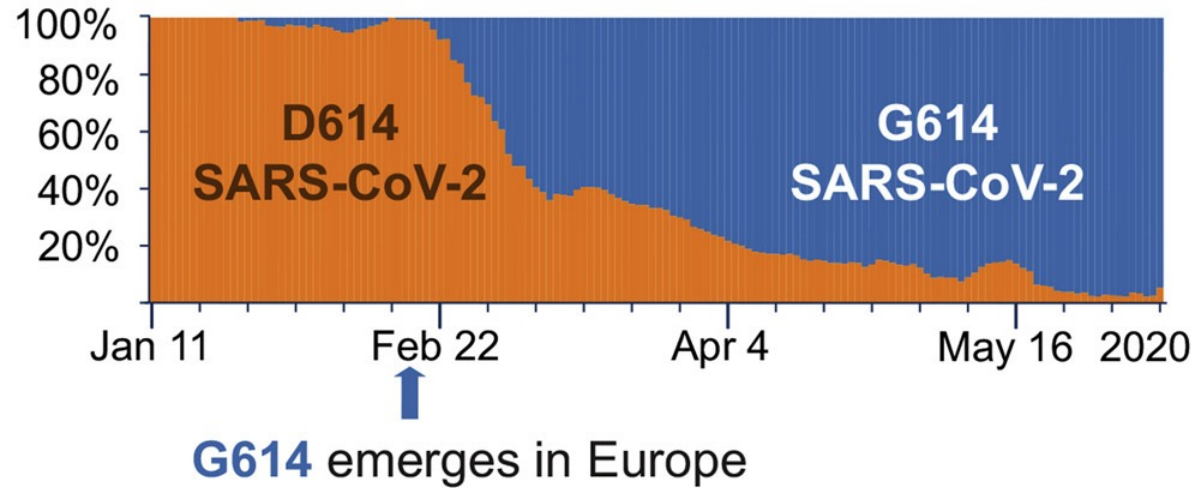
- Pro: parameter estimates map to biologically meaningful quantities (R_0 , length of infectious period, etc)
- Con: computational cost (BDMM) or approximations (MSBD)
- Con: sensitive to unmodelled sampling discrepancies

Phylogenetic tree of SARS-CoV-2 sequences isolated in Washington State reveal multiple introductions.



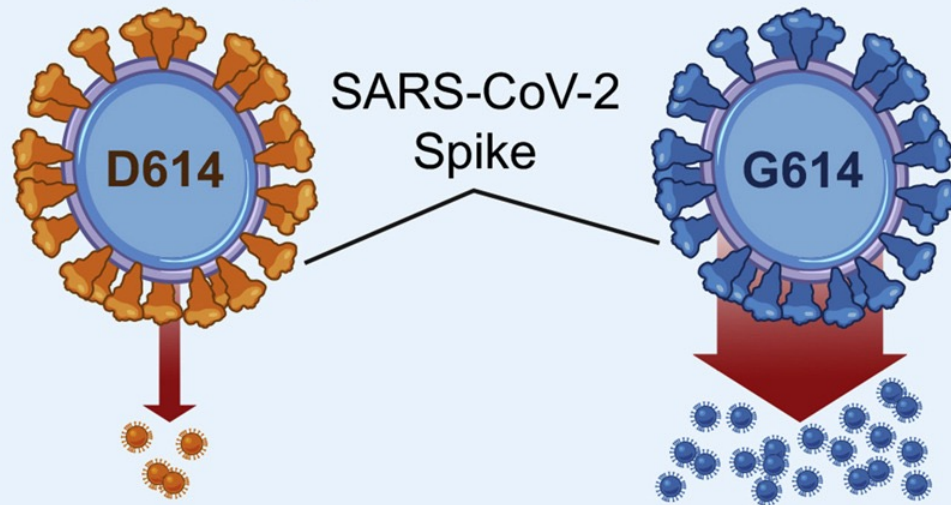
<https://nextstrain.org/groups/blab/ncov/wa-phylodynamics?c=county>

Global Transition



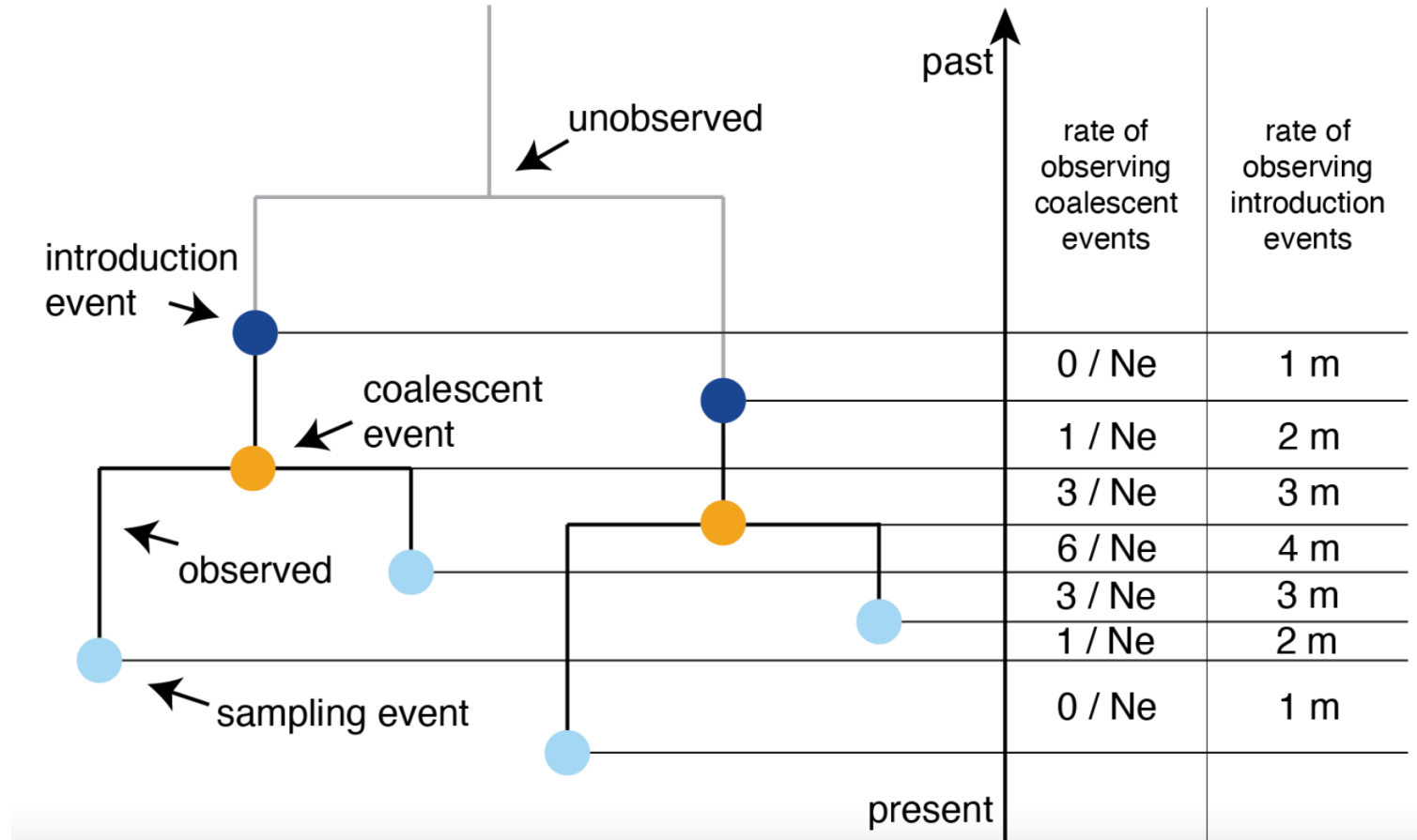
The two spike variants 614D and 614G were shown to differ in viral load and presumably transmission rate.

Magnitude of Infection

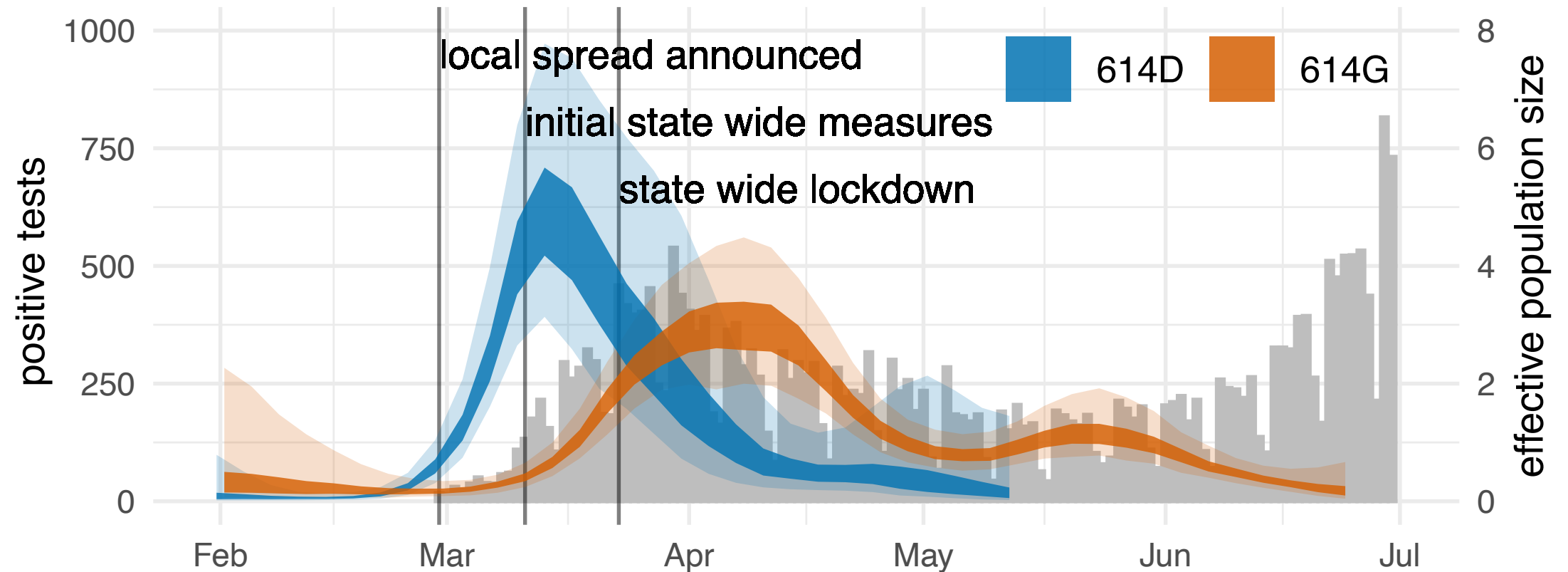


Korber et al. (2020), *Cell*

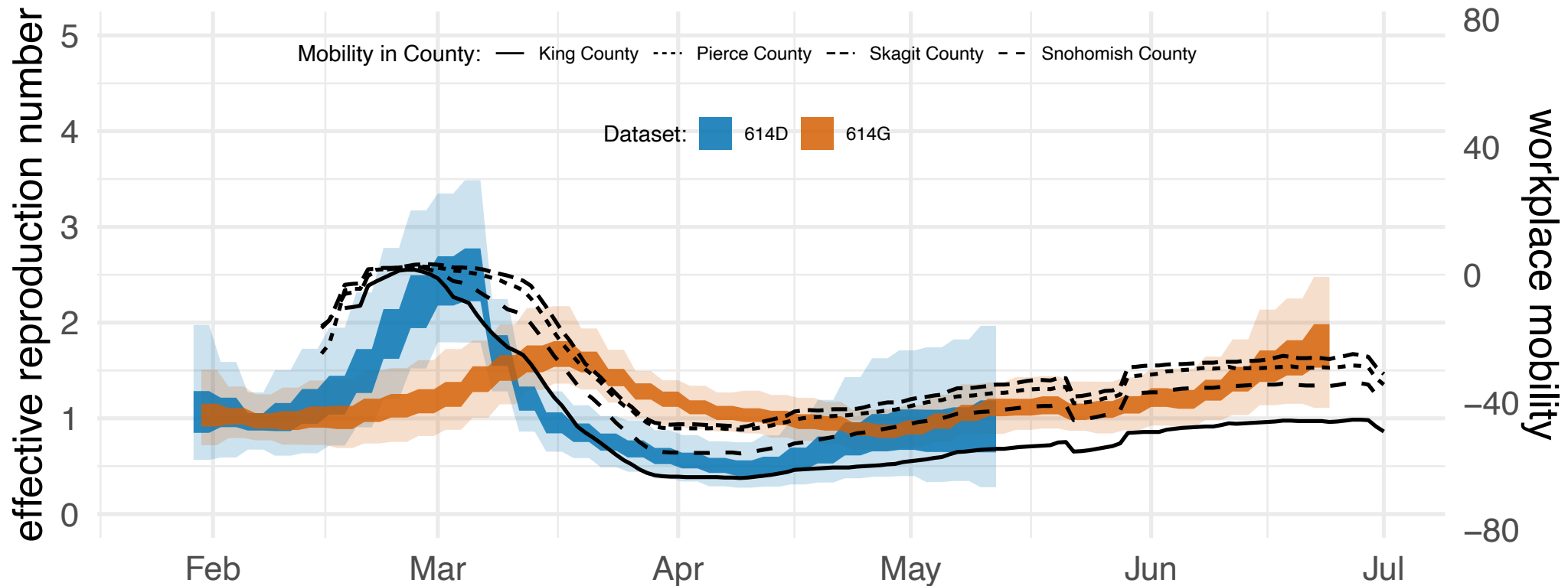
Migration histories can be conditioned on instead of estimated



Reconstructing the effective population size of the spike variants reveals different temporal patterns

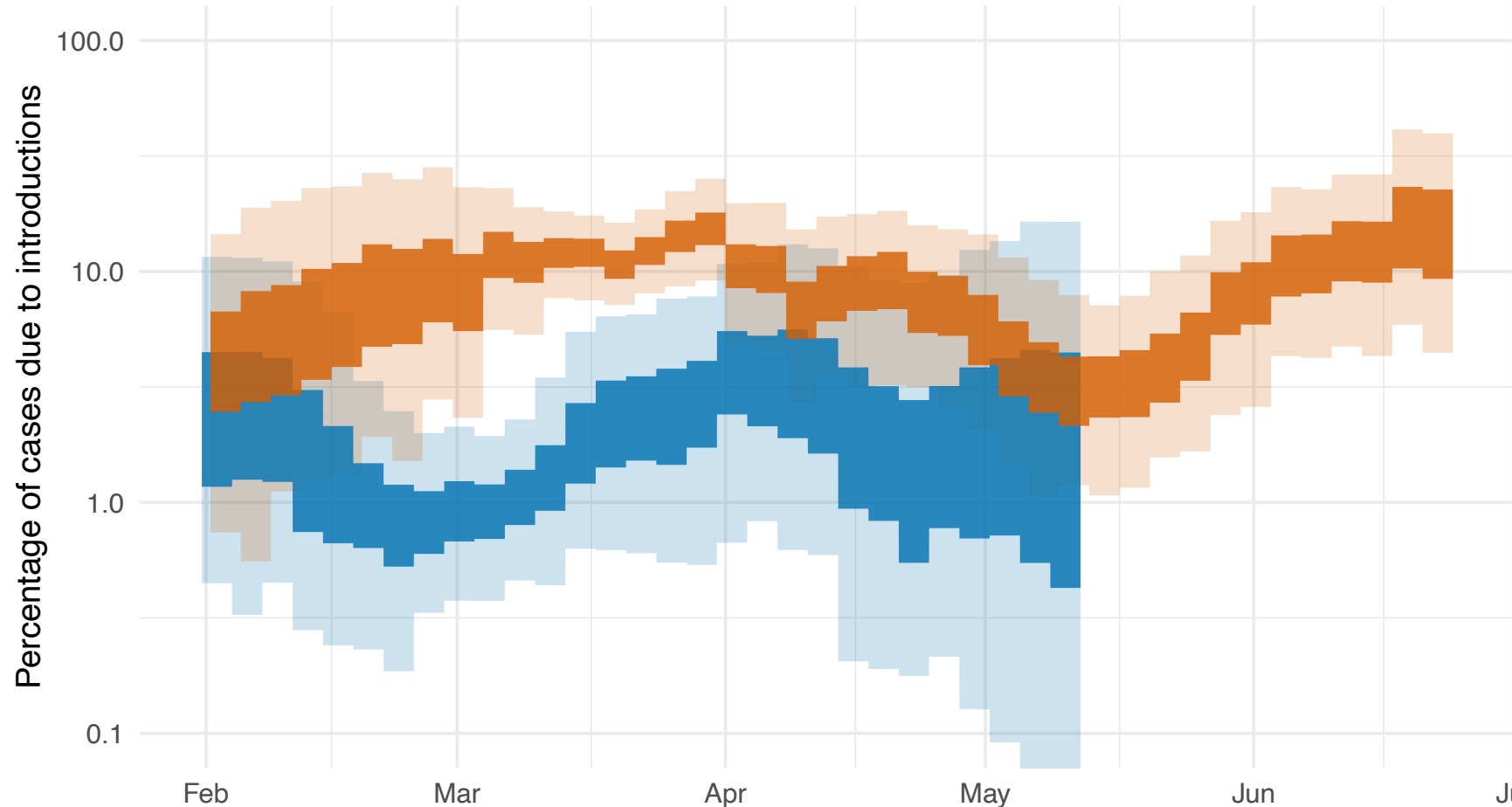


Changes in the reproduction rate of different spike variants reflect different mobility trends in the state



Müller, Wagner, Frazar and Roychoudhury et al. (2021), *Sci. Trans. Med.*

Introductions played different roles in the spread of different SARS-CoV-2 spike variants

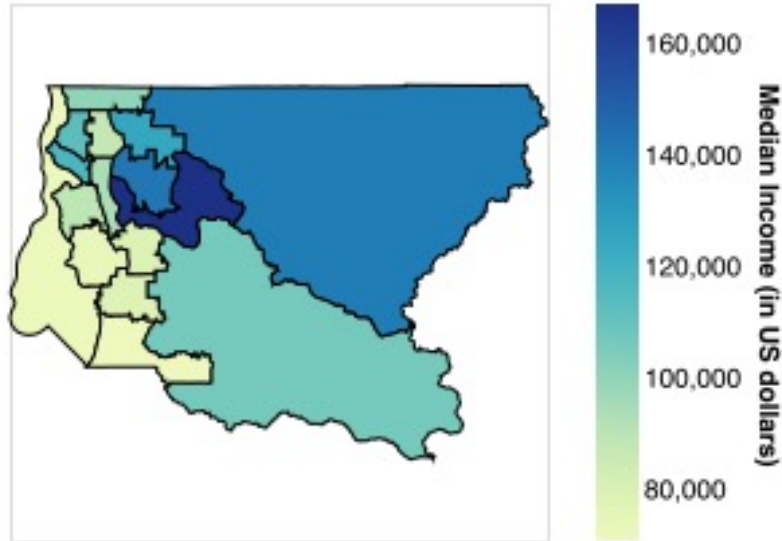


Müller, Wagner, Frazar and Roychoudhury et al. (2021), *Sci. Trans. Med.*

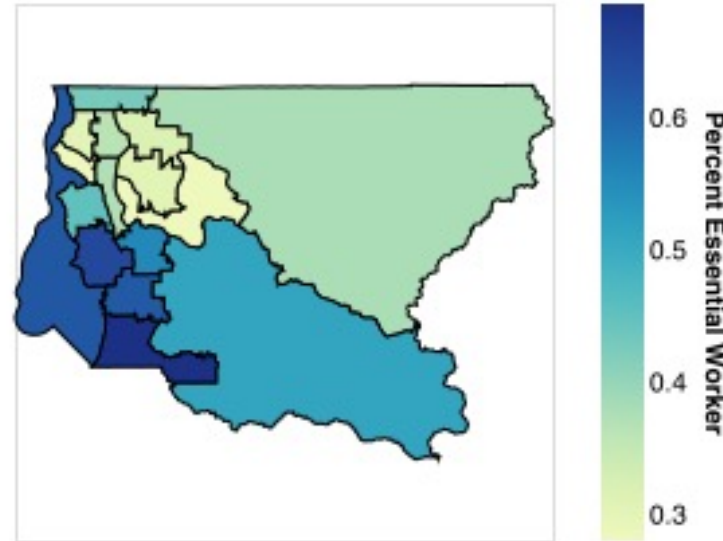
Higher incomes in North King County, fewer essential workers, smaller households (work led by Miguel Paredes)



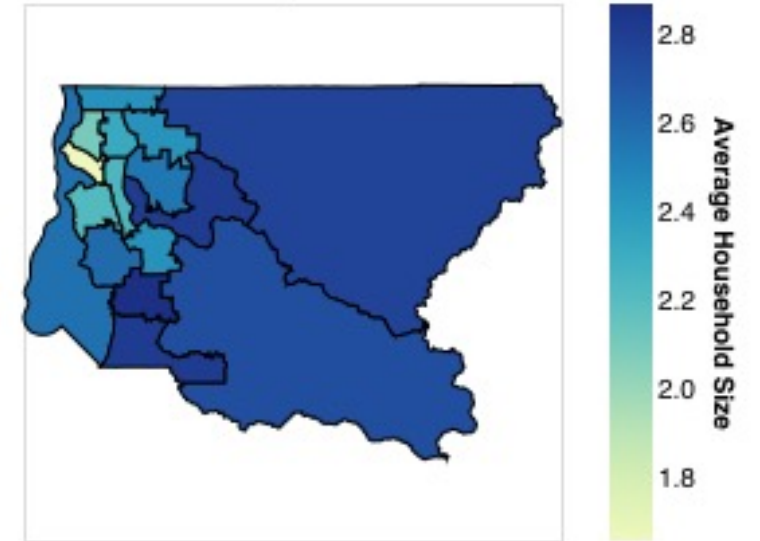
2020 Median Household Income



Percent essential worker 2015-2020

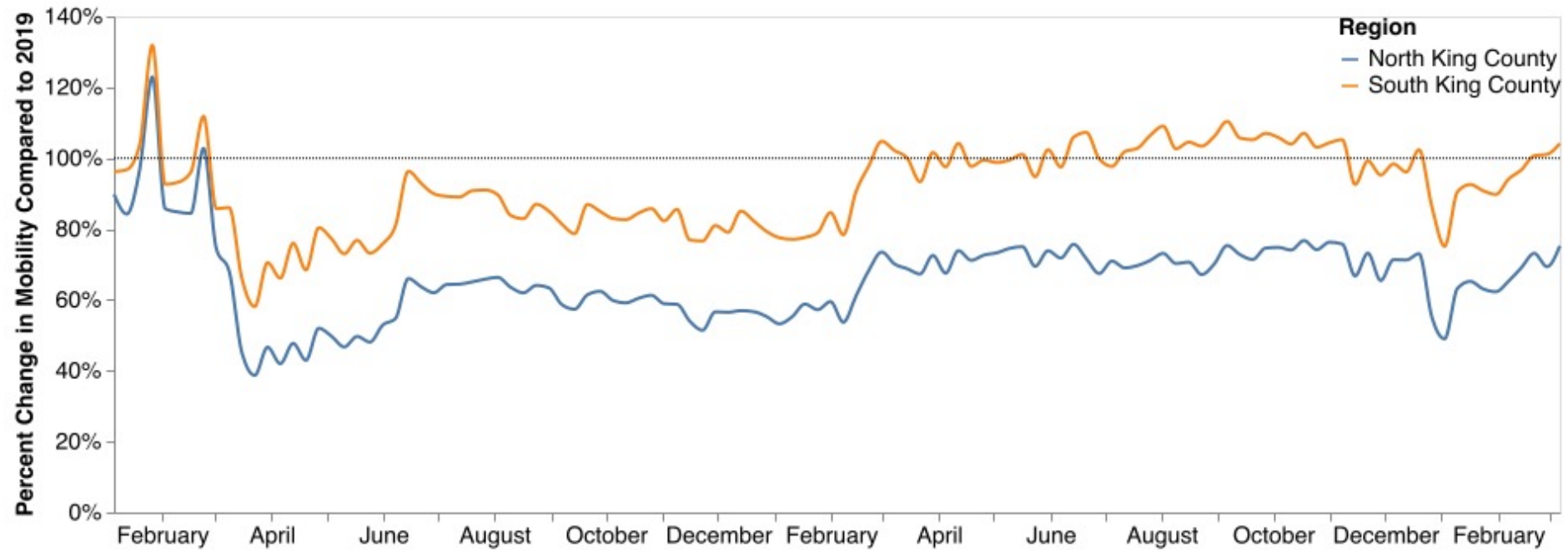


2020-2015 Average Household Size



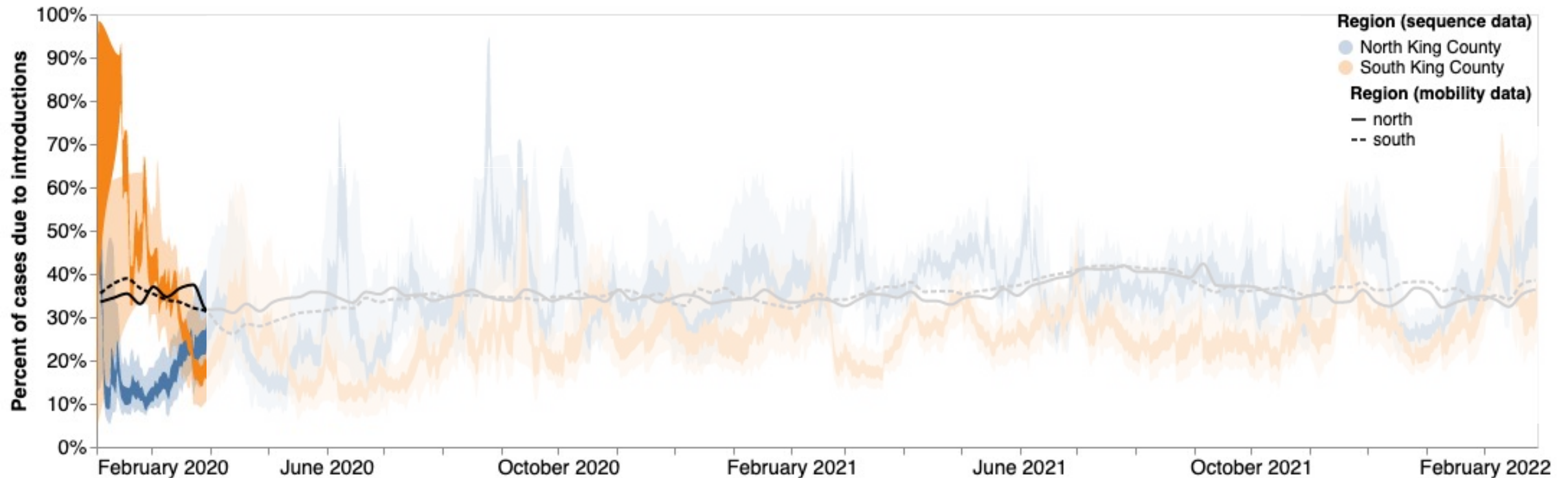
M. I. Paredes, ..., **N. F. Müller***, T. Bedford*, “*Local-Scale phylodynamics reveal differential community impact of SARS-CoV-2 in metropolitan US county*” medRxiv. 2022
*joint supervision

Stronger and more sustained reductions in mobility in North King County



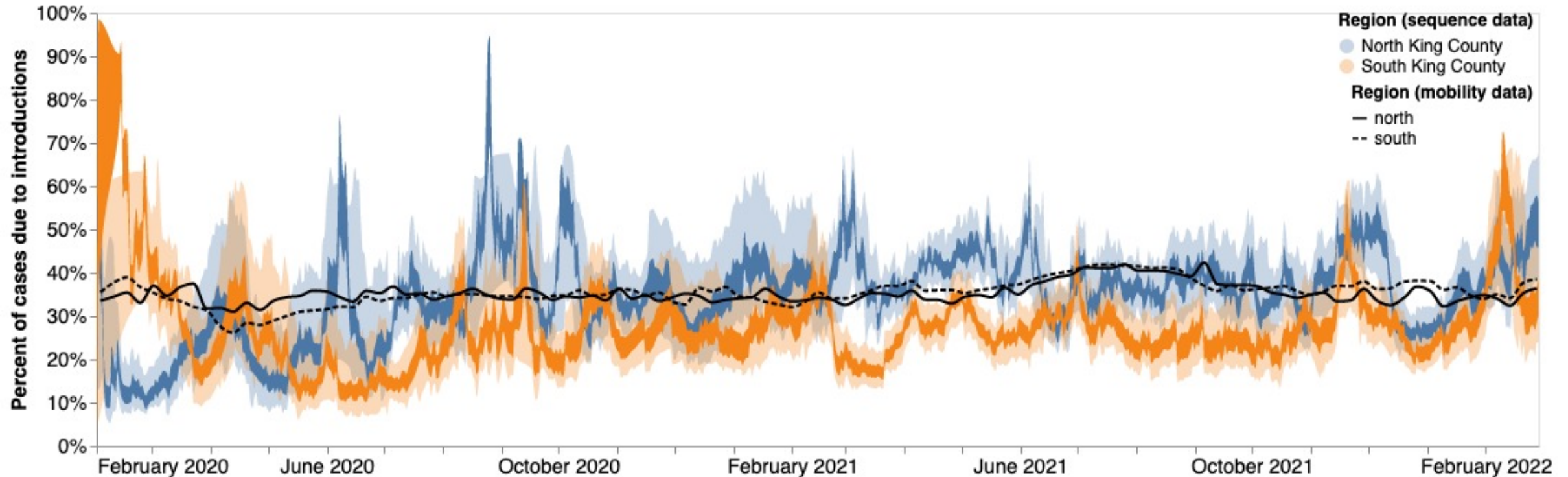
M. I. Paredes, ..., **N. F. Müller***, T. Bedford*, “*Local-Scale phylodynamics reveal differential community impact of SARS-CoV-2 in metropolitan US county*” medRxiv. 2022
*joint supervision

Pre measures, local spread dominated in North, while introductions dominated in South



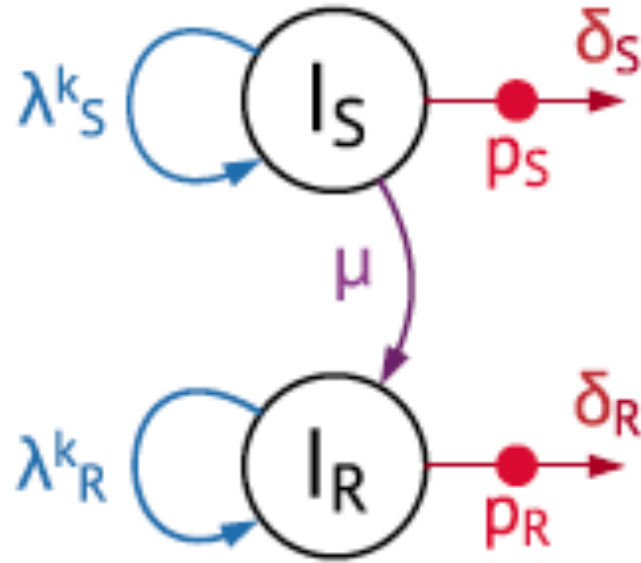
M. I. Paredes, ..., **N. F. Müller***, T. Bedford*, “*Local-Scale phylodynamics reveal differential community impact of SARS-CoV-2 in metropolitan US county*” medRxiv. 2022
*joint supervision

After measures around 30-40% of cases in North were from introductions, while only 20-30% of cases in South were



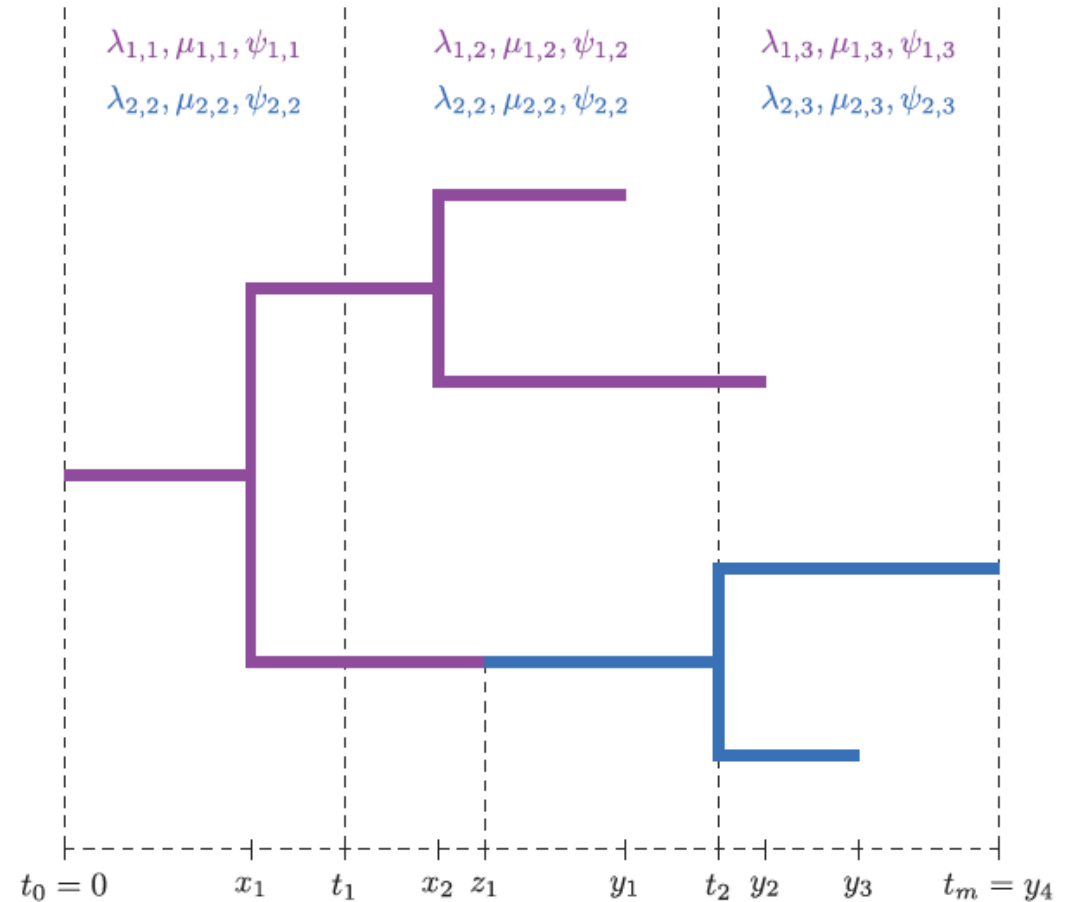
M. I. Paredes, ..., **N. F. Müller***, T. Bedford*, "*Local-Scale phylodynamics reveal differential community impact of SARS-CoV-2 in metropolitan US county*" medRxiv. 2022
*joint supervision

The multi-type birth-death process models how lineages give birth, die, are sampled or jump between states.



Multi-type birth-death process

- N types, $i = 1, \dots, N$
- N birth rates λ_i
- N death rates μ_i
- N sampling rates ψ_i
- Migration rates $m_{i,j} \forall i \neq j$



QUESTIONS?

Some reading material

- Paper on DTA: <https://doi.org/10.1371/journal.pcbi.1000520>
- How do different structured coalescent methods compare: <https://academic.oup.com/mbe/article-abstract/34/11/2970/3896419>
- Birth-death models with migration: <https://doi.org/10.1093/molbev/msw064>
- GLM for DTA: <https://doi.org/10.1371/journal.ppat.1003932>
- GLM for MASCOT: <https://doi.org/10.1093/ve/vez030>