



UNIVERSITY OF
OXFORD



BIG DATA
INSTITUTE



Analysing within- and between-host HIV genetic diversity with phyloscanner for detection of transmission and more

Chris Wymant^{1,2} on behalf the BEEHIVE and PANGEA Collaborations

¹ Big Data Institute, Li Ka Shing Centre for Health Information and Discovery,
Nuffield Department of Medicine, University of Oxford

² Department of Infectious Disease Epidemiology, Imperial College London

The BEEHIVE Project: Bridging the Epidemiology and Evolution of HIV in Europe

Oxford University

Christophe Fraser
Tanya Golubchik
Matthew Hall
Michelle Kendall
Rob Power
Chris Wymant

Imperial College London

Paul Kellam
Frank de Wolf

Amsterdam Medical Centre

Margreet Bakker
Ben Berkhout
Marion Cornelissen
Peter Reiss

Wellcome Trust Sanger Institute

Swee Hoe Ong

European Bioinformatics Institute

Astrid Gall
Martin Hunt

HIV Monitoring Foundation

Daniela Bezemer
Mariska Hillebregt
Ard van Sighem
Sima Zaheri

Karolinska Institute

Jan Albert

Antwerp Institute of Tropical Medicine

Katrien Fransen
Guido Vanham

John Hopkins University

M. Kate Grabowski

Robert Koch-Institute, Berlin

Norbert Bannert
Claudia Kücherer

University Hospital Zürich

Huldrych Günthard
Roger Kouyos

Division of Intramural Research NIAID, Baltimore

Oliver Laeyendecker

Helsinki University Hospital

Pia Kivelä
Kirsi Liitsola
Matti Ristola

Université Paris Sud

Laurence Meyer

University College London

Kholoud Porter

College de France

François Blanquart

École polytechnique fédérale de Lausanne

Jacques Fellay

Analysis Advisory Group

Samuel Alizon
Sebastian Bonhoeffer
Gabriel Leventhal
Andrew Rambaut
Oliver Pybus
Gil McVean

With thanks to

Nick Croucher
Katrina Lythgoe
Oliver Ratmann

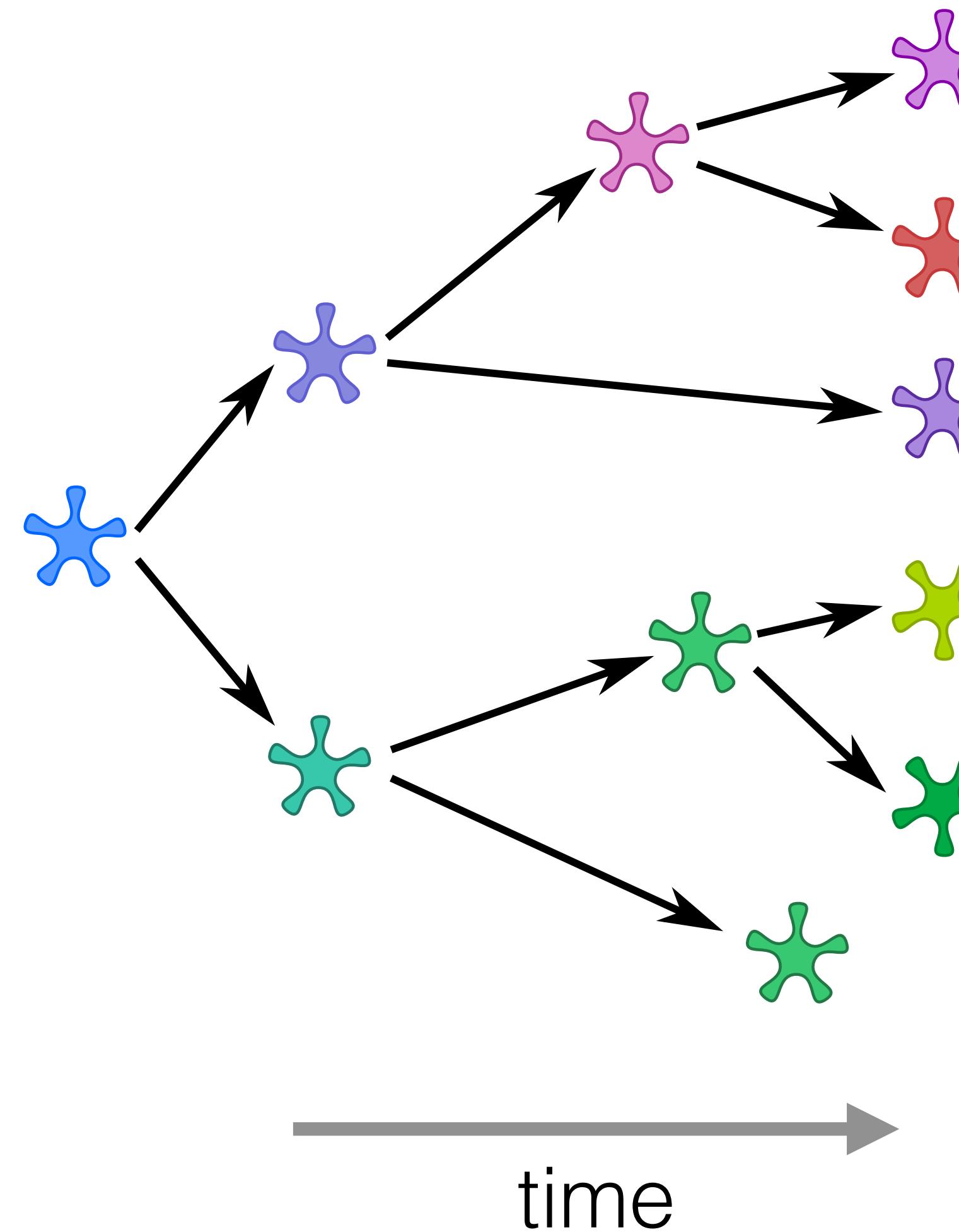


Problem: how can we target our interventions in an epidemic to have a greater impact?

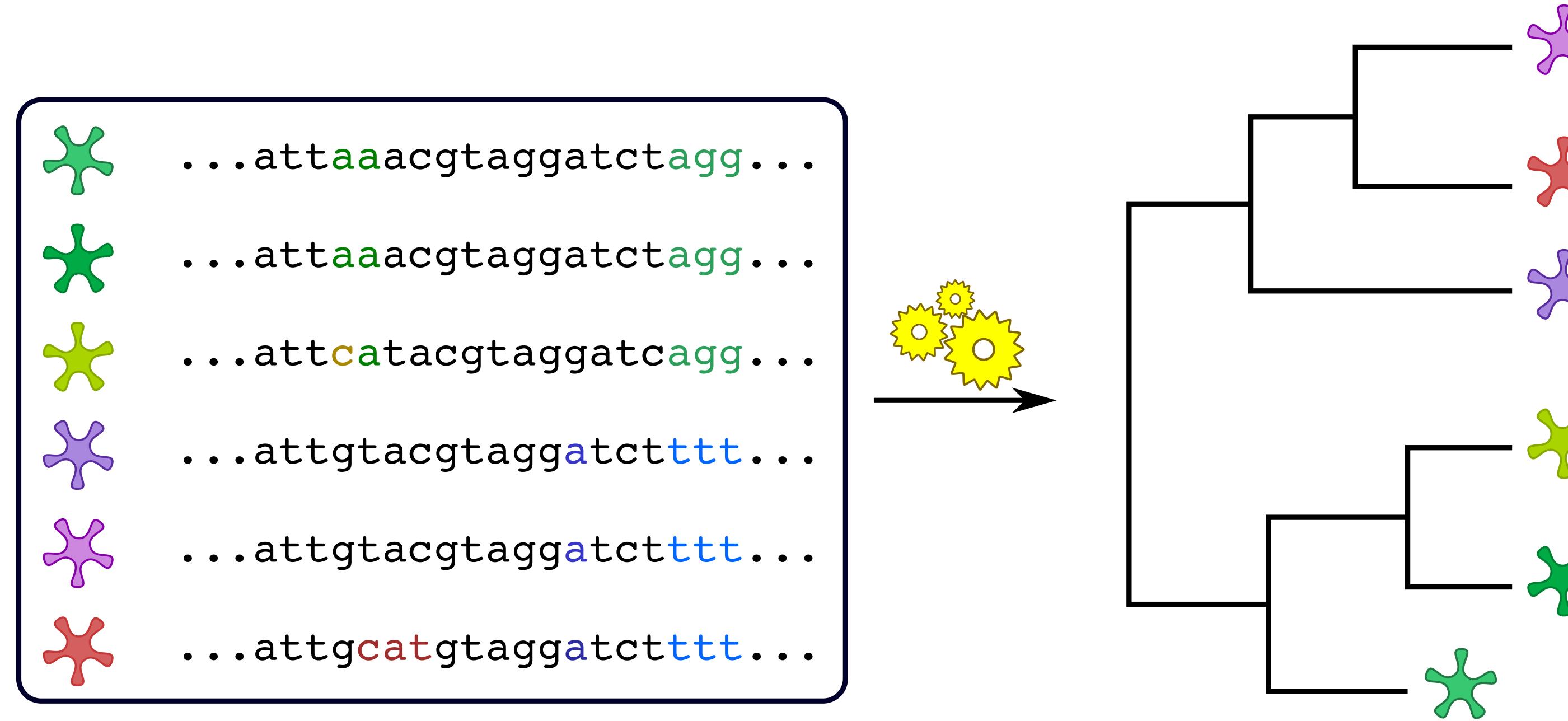
A solution: learn more about the epidemic through analysis of pathogen sequence data.

Genetic changes (substitution) accumulate over time

Substitution: replacement of a nucleotide (e.g. a → t)

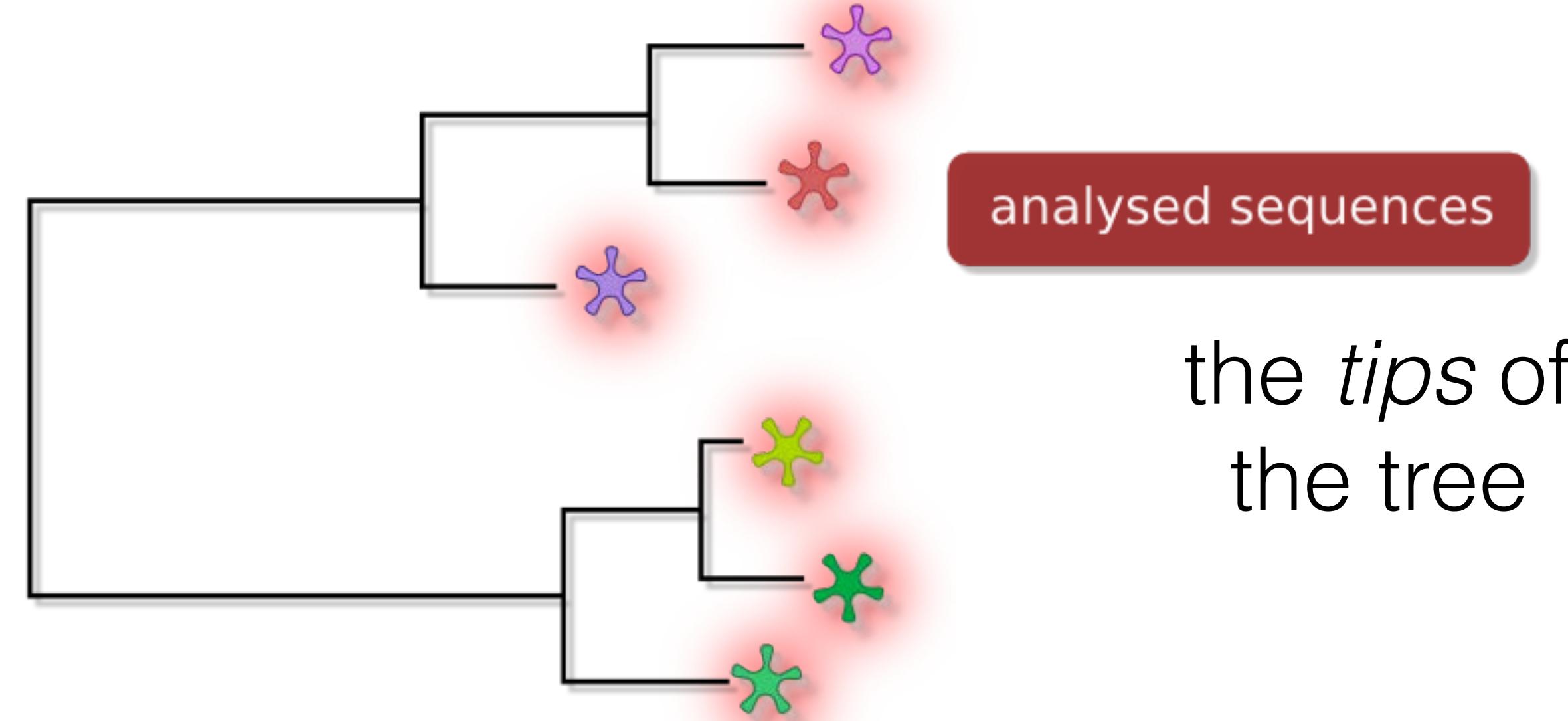


Using substitution patterns to reconstruct the evolutionary history



Phylogenetics aim to reconstruct evolutionary trees (*phylogenies*) from genetic sequence data.

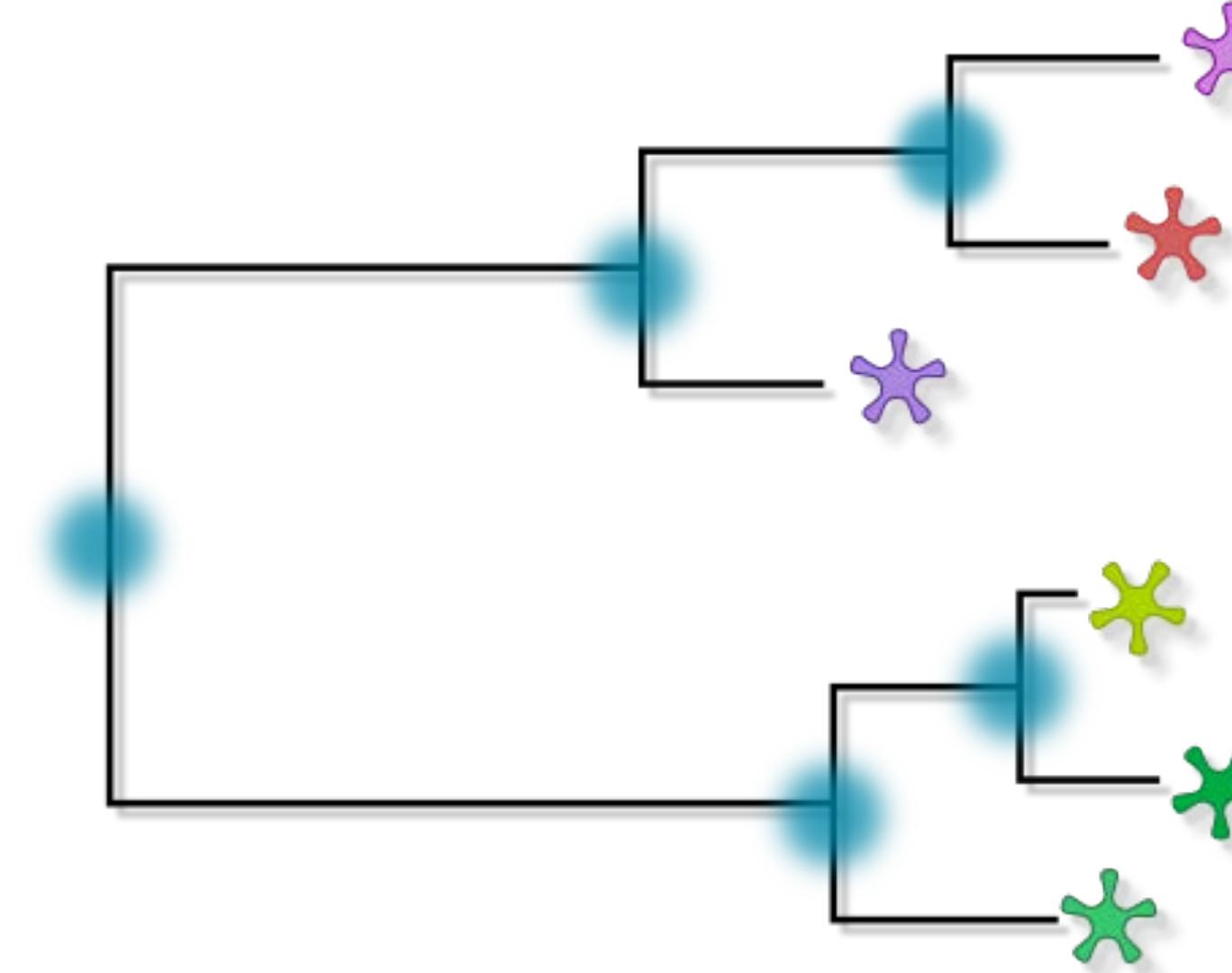
Using trees to represent the evolutionary history



the *tips* of
the tree

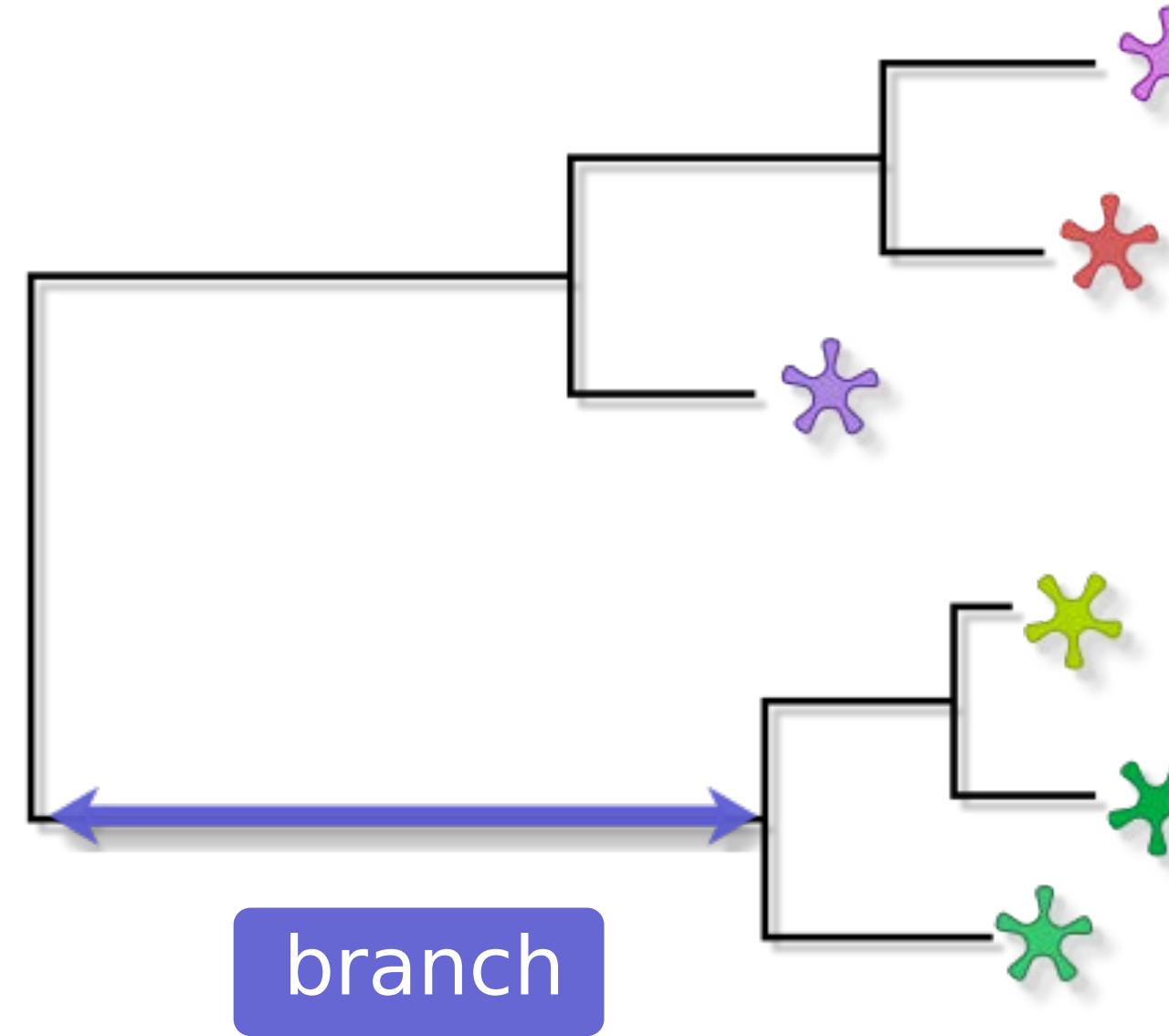
Using trees to represent the evolutionary history

Most Recent Common Ancestors
(MRCA)



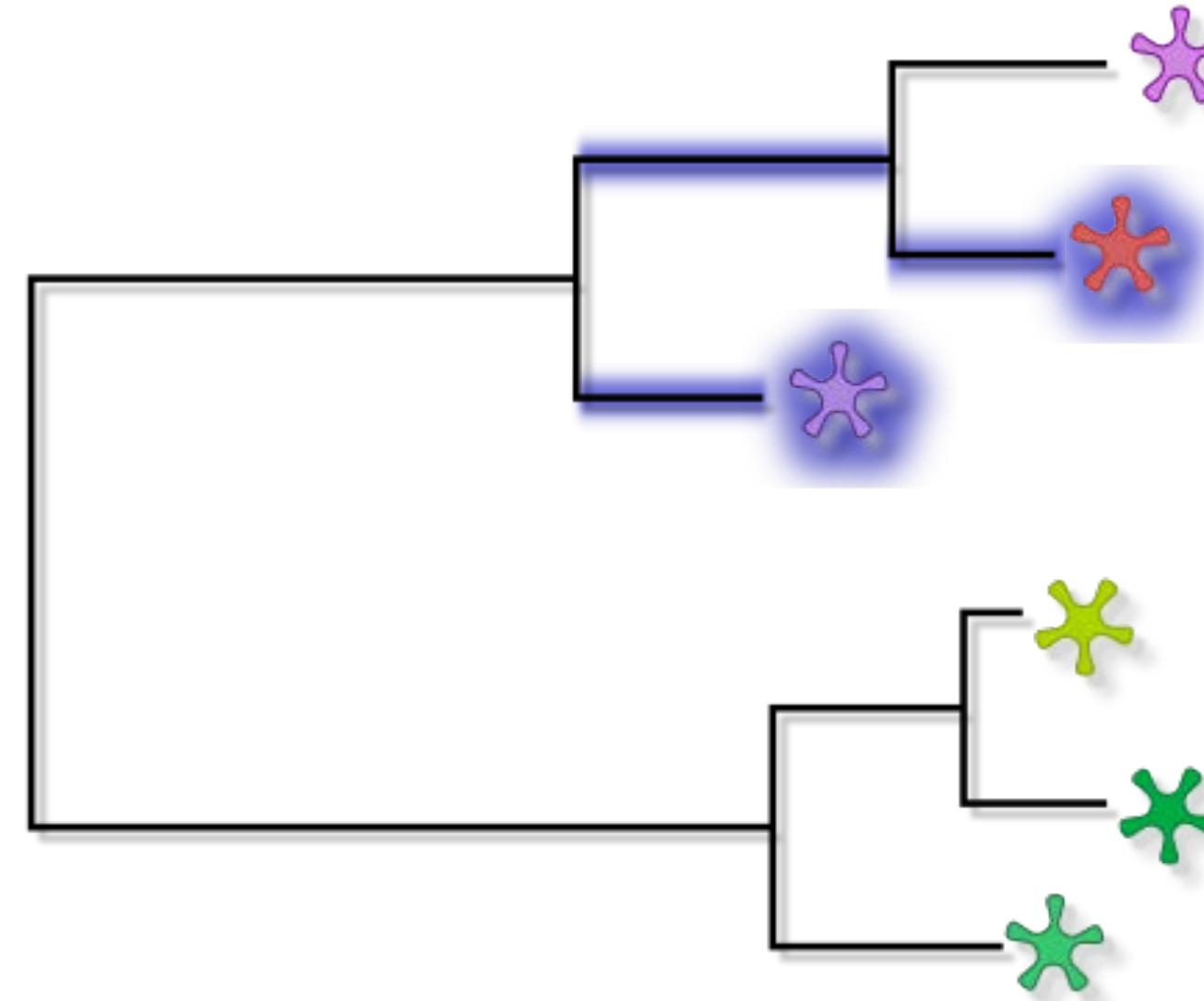
the *nodes*
of the tree

Using trees to represent the evolutionary history



length = amount of evolution (**not time**, as a rule)

Using trees to represent the evolutionary history



distances between tips

"*patristic*" distance: sum of branch lengths

Original slide by
Thibaut Jombart

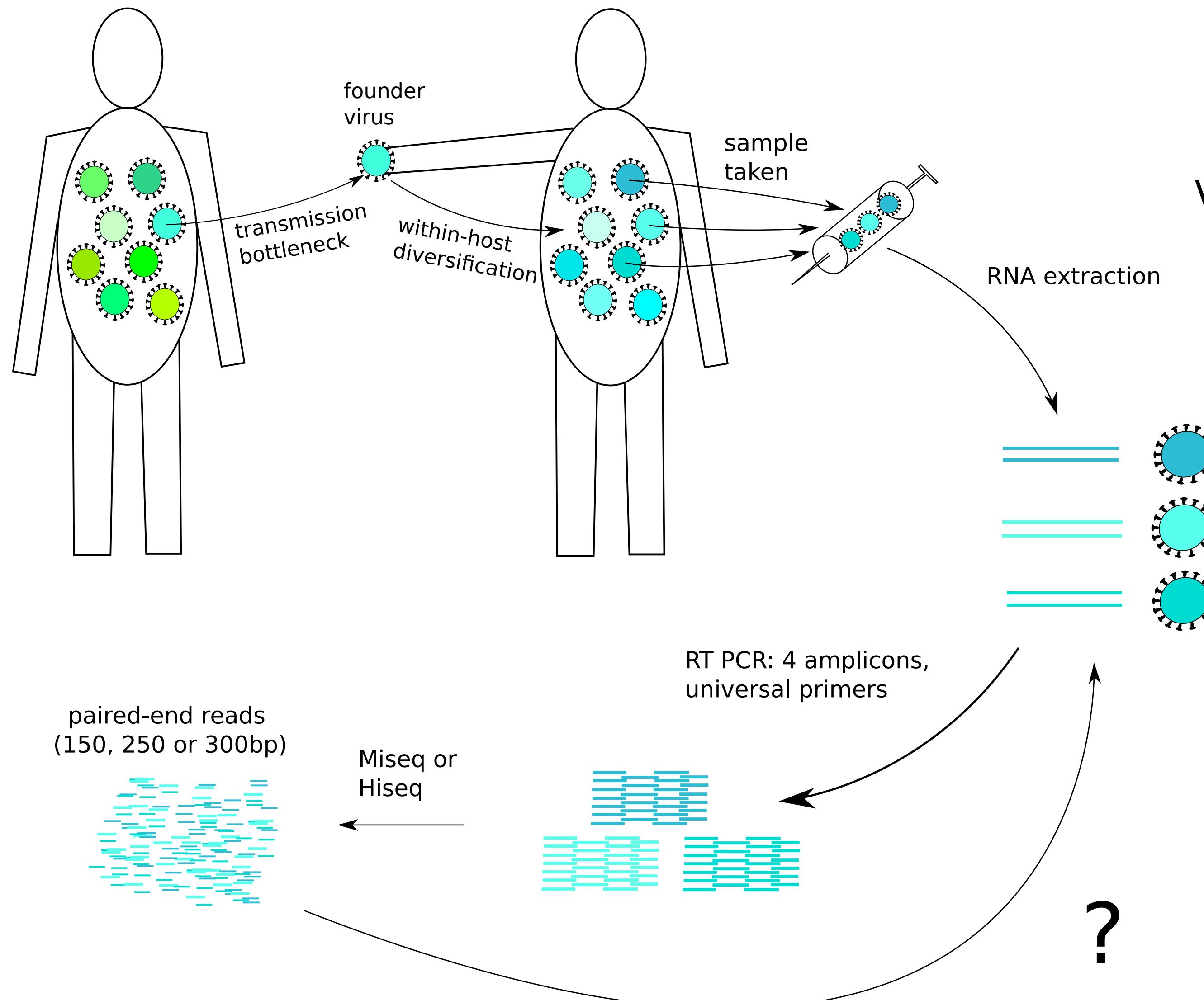
Small distance \Rightarrow small amount of evolution,
from which we infer proximity in a transmission network.

Problem: to meaningfully interpret small differences between closely related viruses, we want *accurate* genomes.

Problem: to make robust statements, need *many* genomes.

Solution: high-throughput sequencing + accurate, scalable sequence processing.

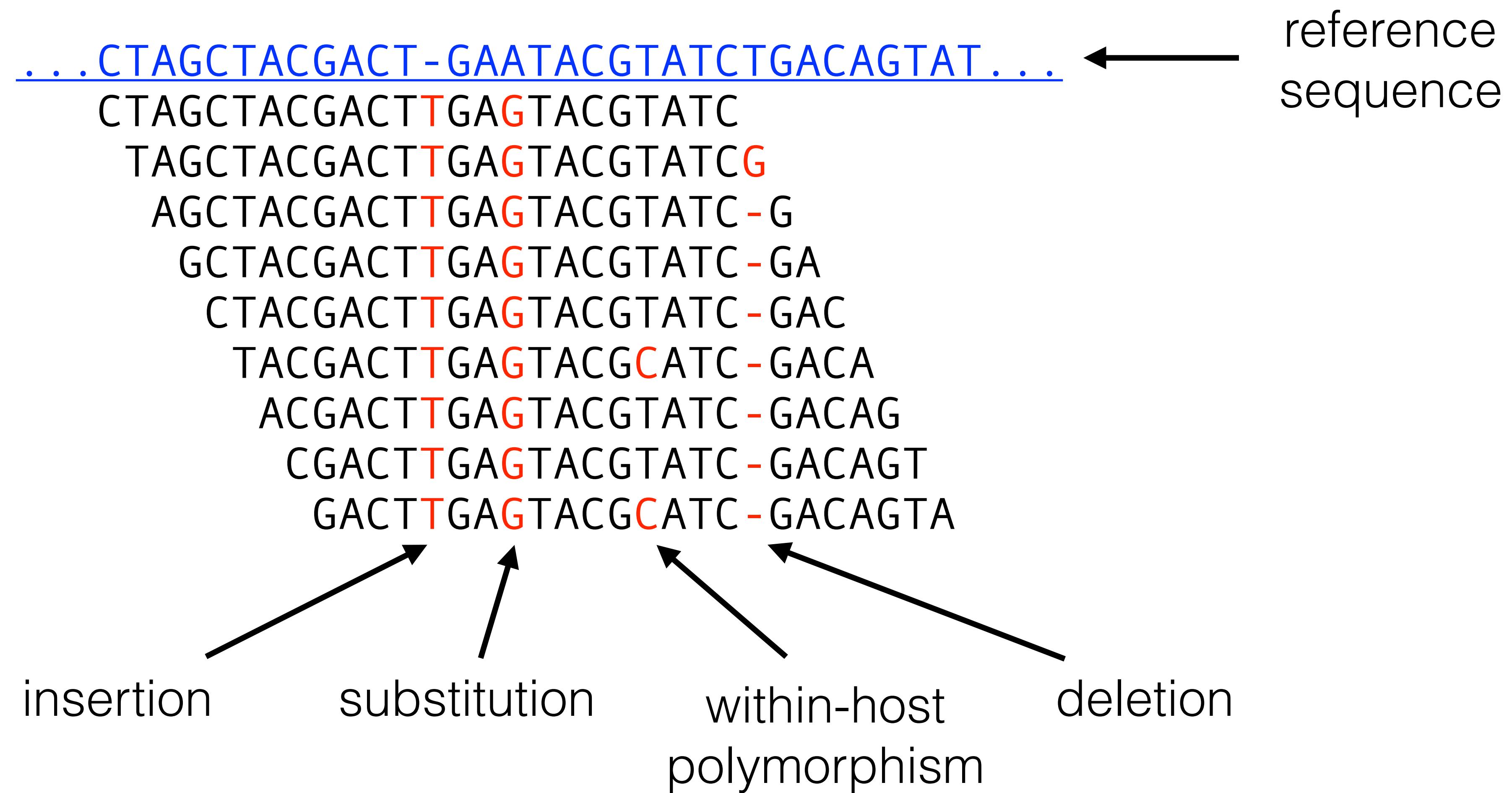
Our HIV Data



Cornelissen *et al.*
Virus Research 2016

Gall *et al.*
J. Clin. Microbiol.
2012

Mapping Reads to a Reference

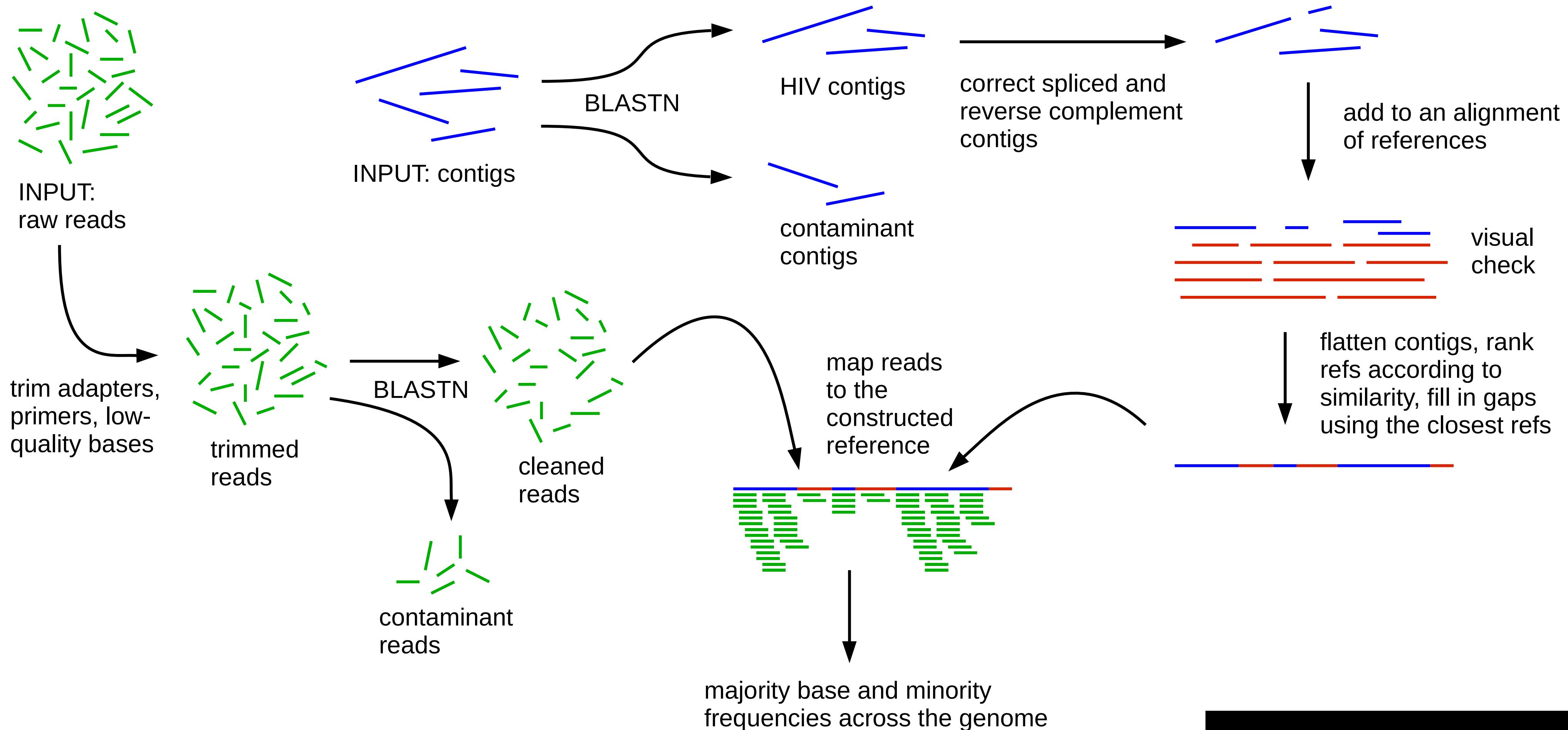


However, the more different a read is from its reference, the more likely it is to be aligned incorrectly or not at all. So **mapping causes biased loss of information.**

e.g. use the same reference for mapping for a group of individuals: bias their viruses to look similar to each other
→ false inference of transmission.

e.g. use old references for mapping new samples
→ false inference of slow evolution

shiver - Sequences from HIV Easily Reconstructed



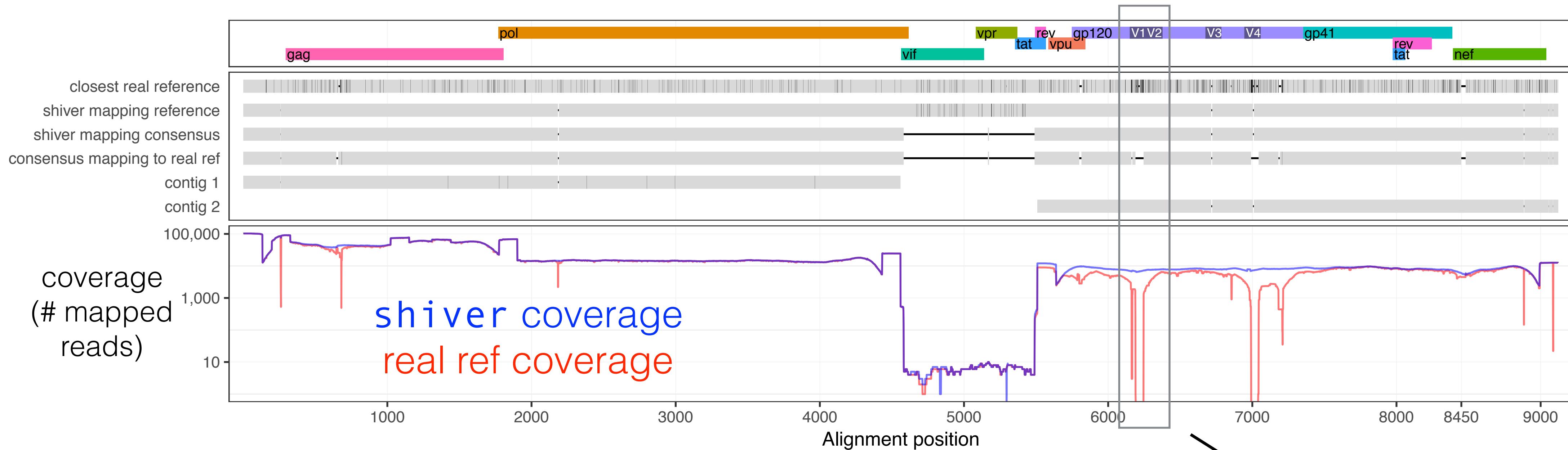
Wymant *et al.*, Virus Evolution 2018

github.com/ChrisHIV/shiver

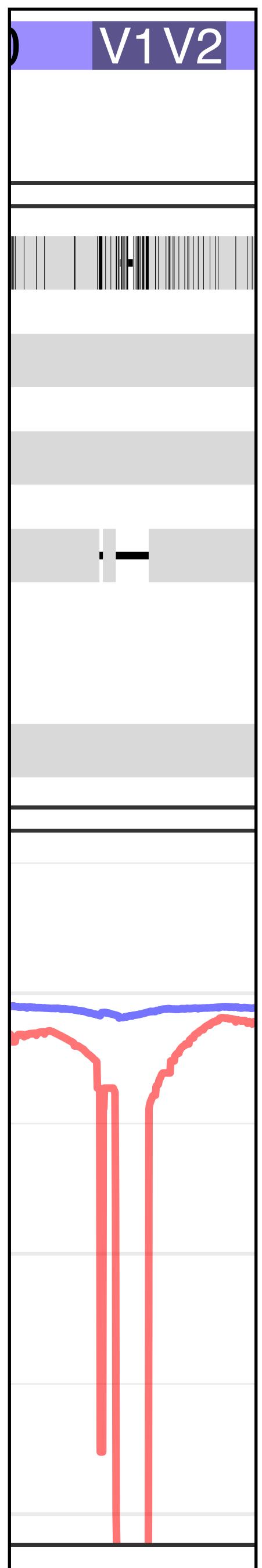
Also works on HCV, RSV, ...

```
$ shiver_align_contigs.sh  
$ shiver_map_reads.sh
```

Map to shiver's constructed reference,
also map to the closest of 3259 real references (from database),
compare.



Over 50 BEEHIVE samples + 65 public samples:
median number of bases called differently with higher coverage:
13; with lower coverage, 0.
Recover missing sequence, often in envelope (vaccine design!).



Application of shiver for 3 HIV projects:

Investigate the viral-genetic basis of virulence

BEEHIVE: $N \sim 3,000$

PopART: $N \sim 1,000$ currently, 9,000 eventually

PANGEA: $N \sim 13,000$ currently, 23,000 eventually

Help interpret the effect of a population-level test-and-treat intervention in Zambia and South Africa

Generate new insights into transmission dynamics in generalized epidemics in Africa



Tanya
Golubchik

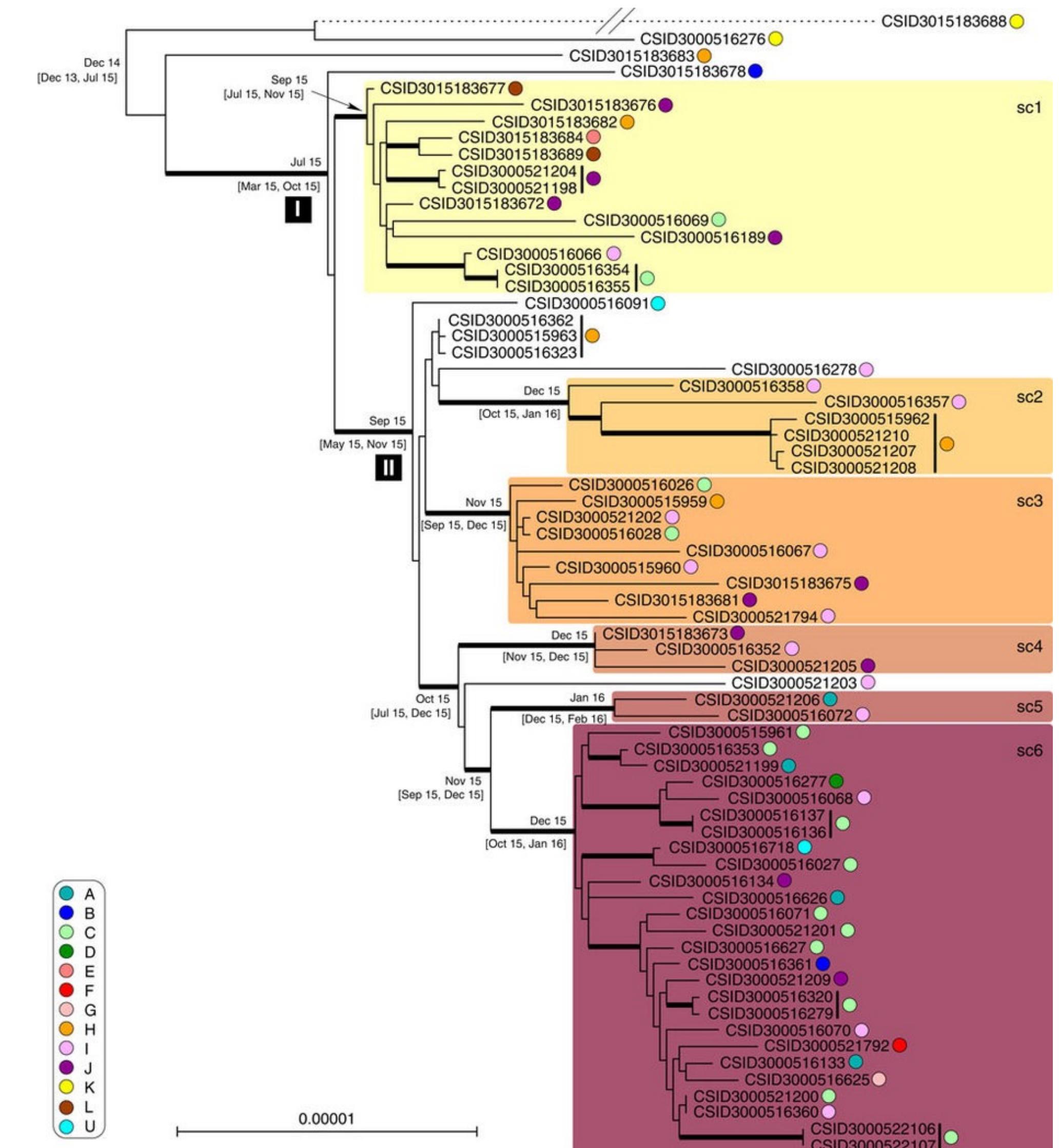
Problem: we want to identify risk factors for transmission. We need to identify not just transmission pairs or clusters, but *who infected whom*.

Molecular epi with one pathogen sequence per individual

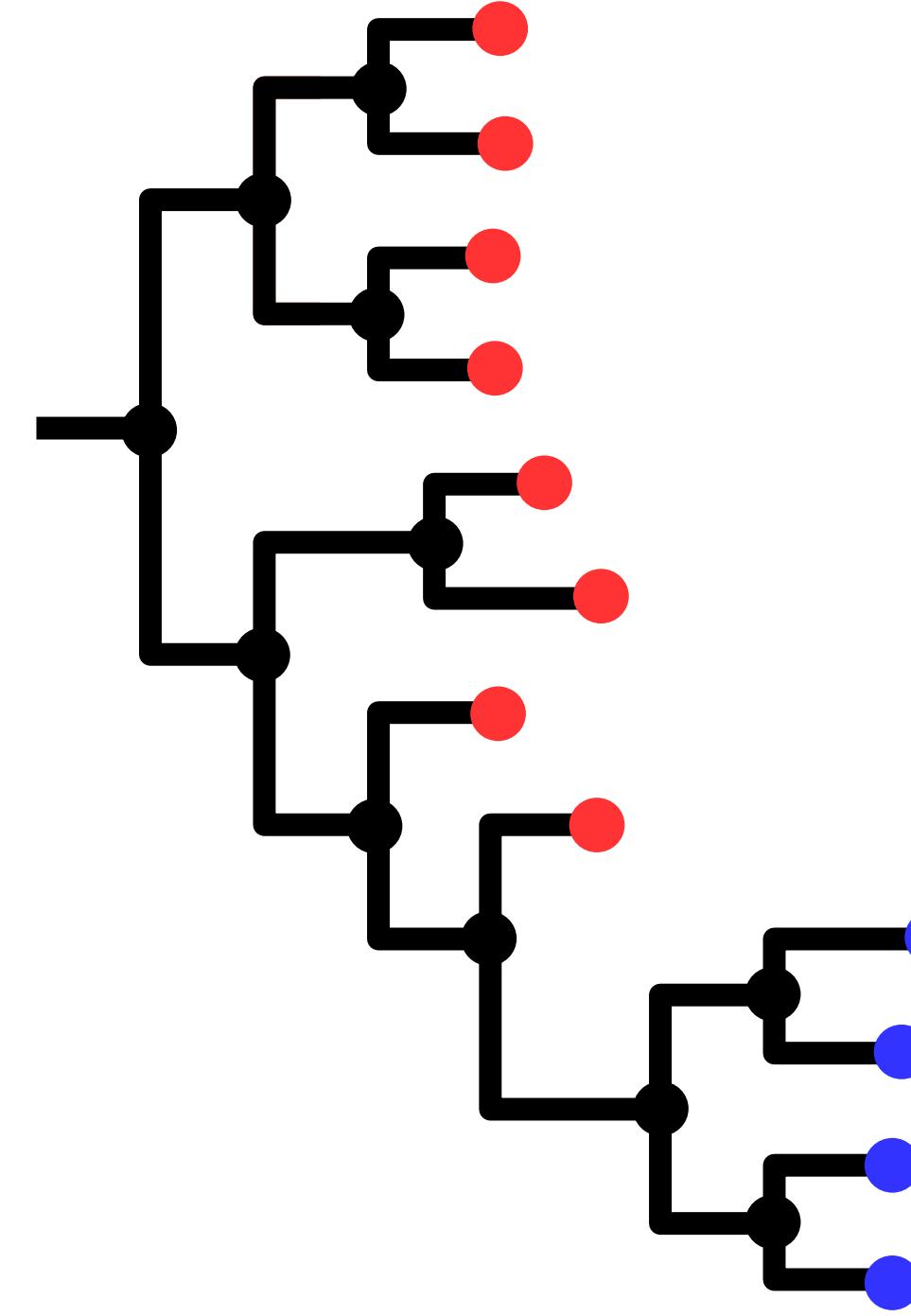
- Make a tree using one pathogen sequence per sampled individual.
- Define clusters: likely to be related by recent transmission.
- Investigate epidemiological correlates.

No information on who is transmitting!

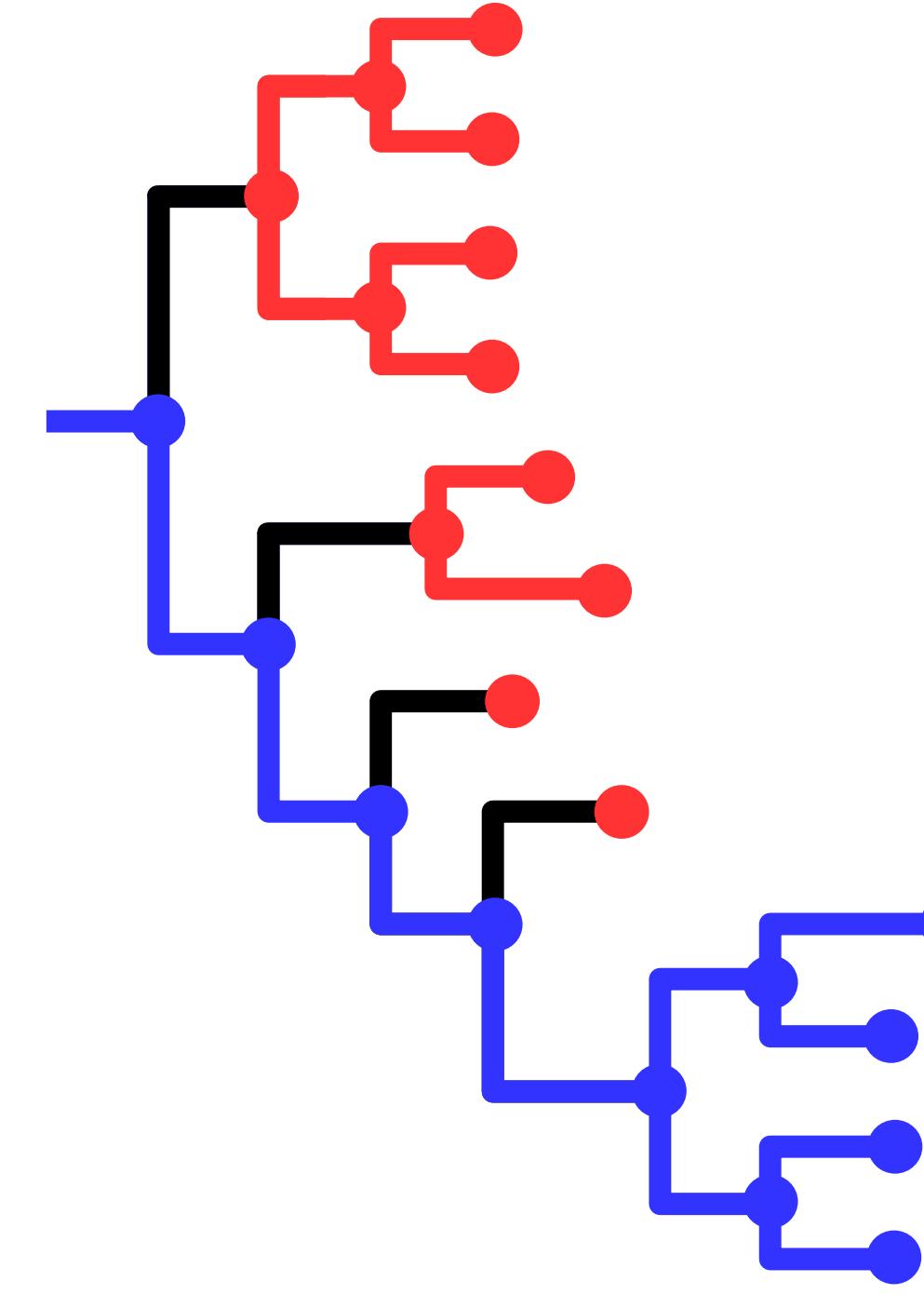
Can do better by also using extra data + transmission model. Conclusions then dependent on both.



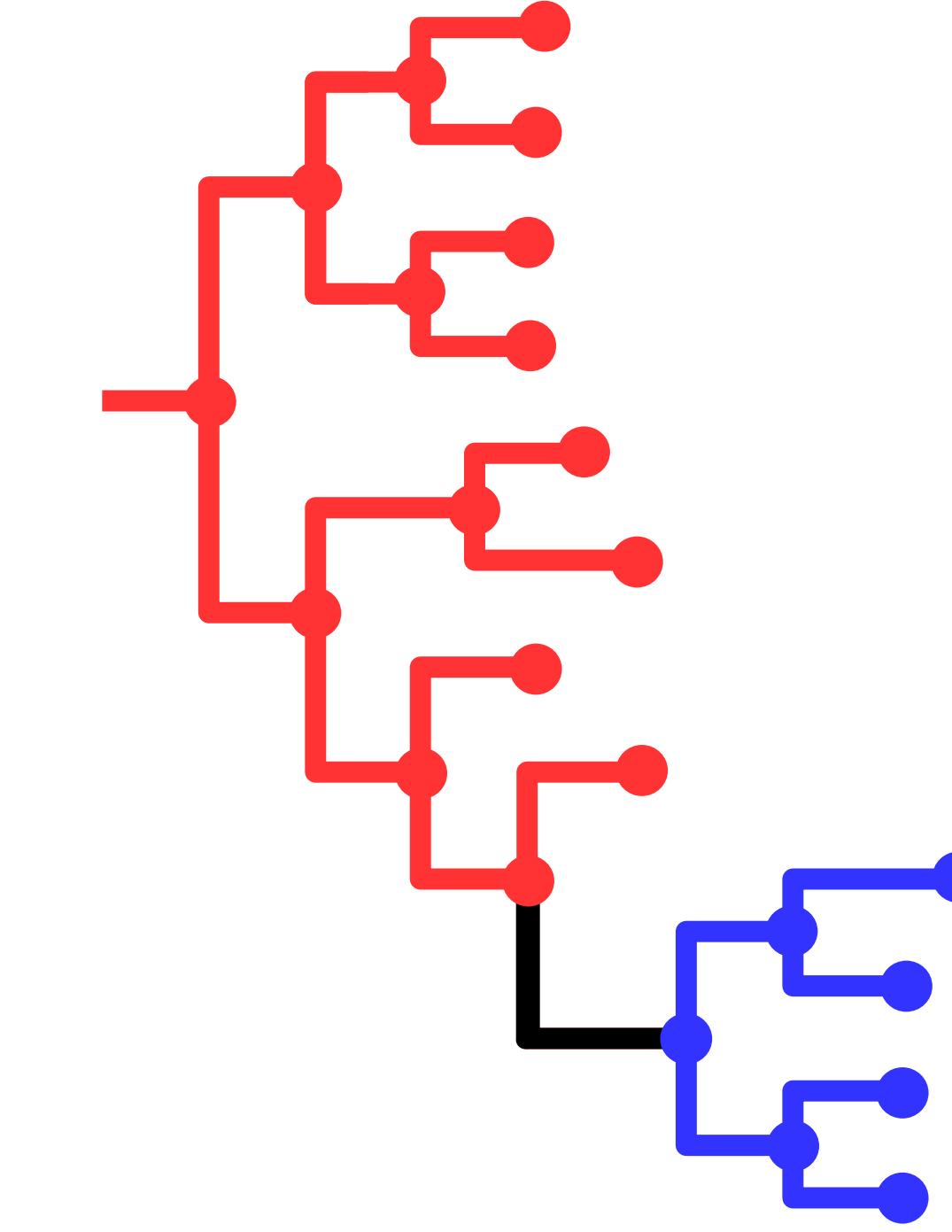
Ancestral State Reconstruction



Known states (red or blue) at the tips.



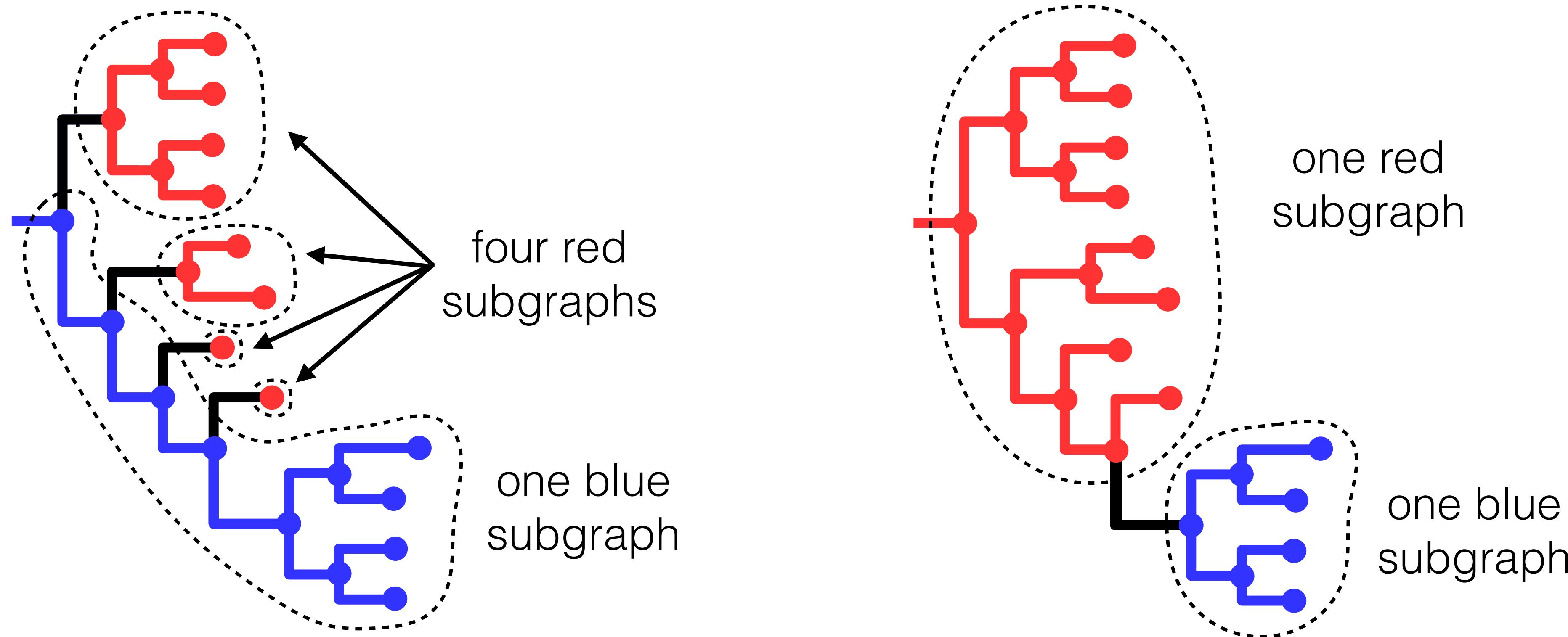
Red evolved from blue: 4+ changes of state needed.



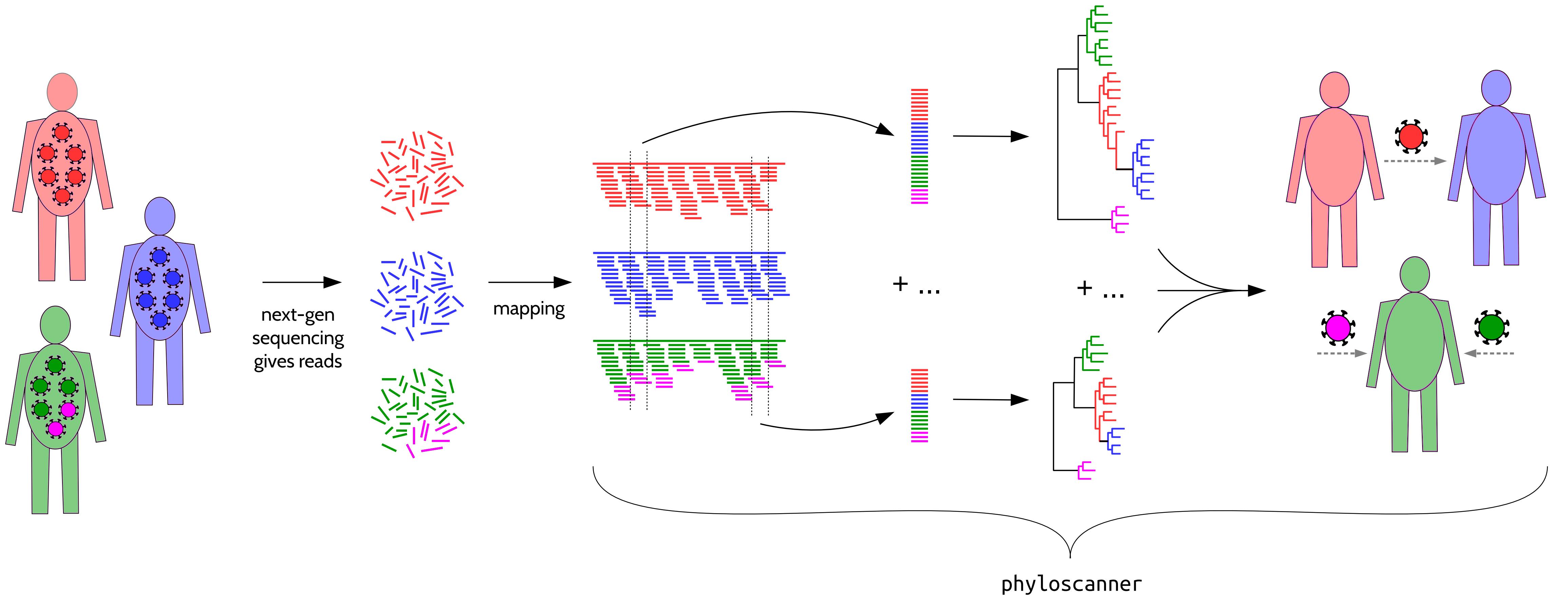
Blue evolved from red: 1 change of state needed.

e.g. colour for “which person was this virus in?” - Romero-Severson et al. PNAS 2016

Ancestral State Reconstruction: subgraphs



A *subgraph* is connected/continuous region of the phylogeny. Later, we'll define these regions by the state assigned to them, i.e. one subgraph = one solid block of colour.

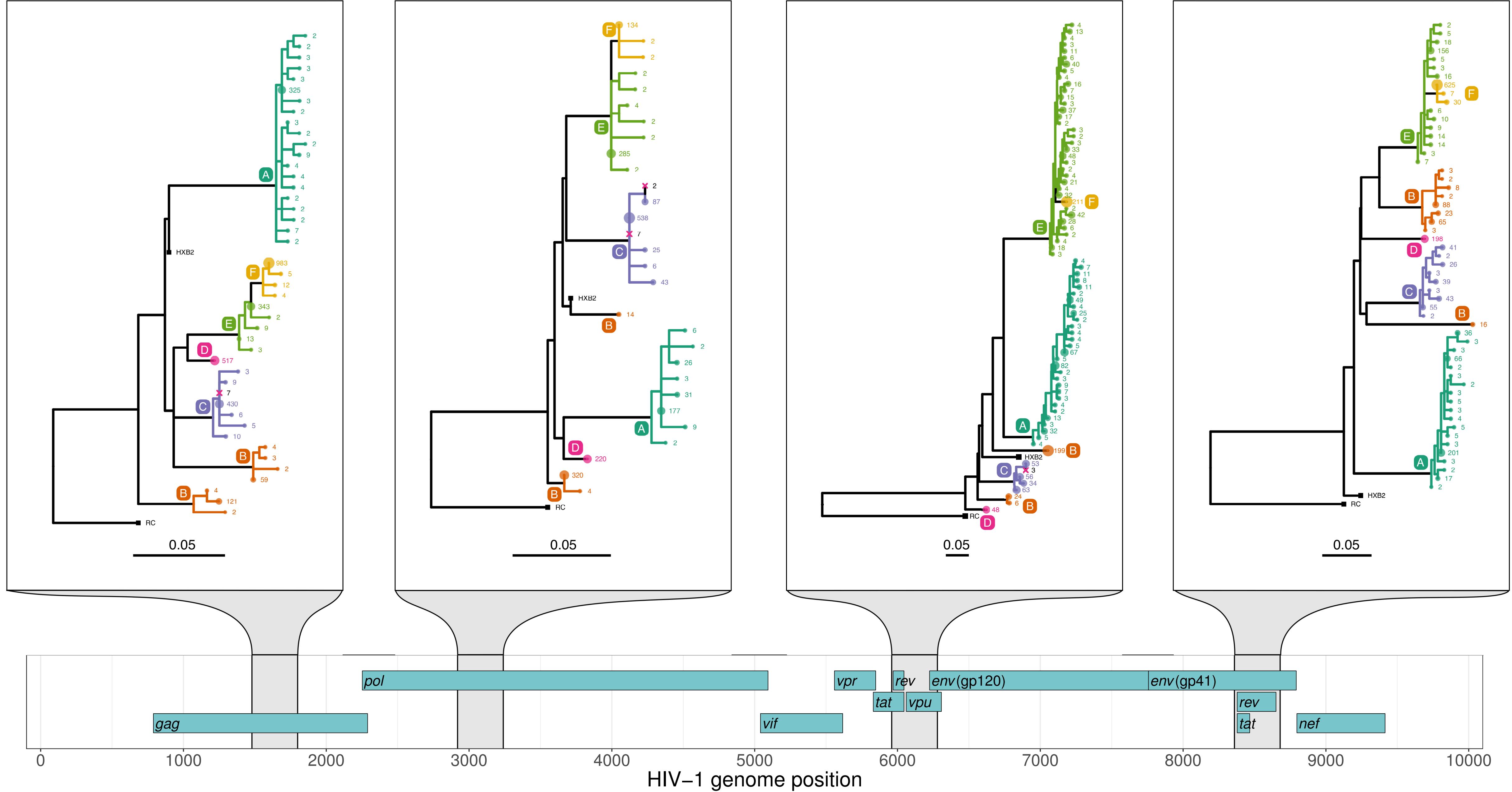


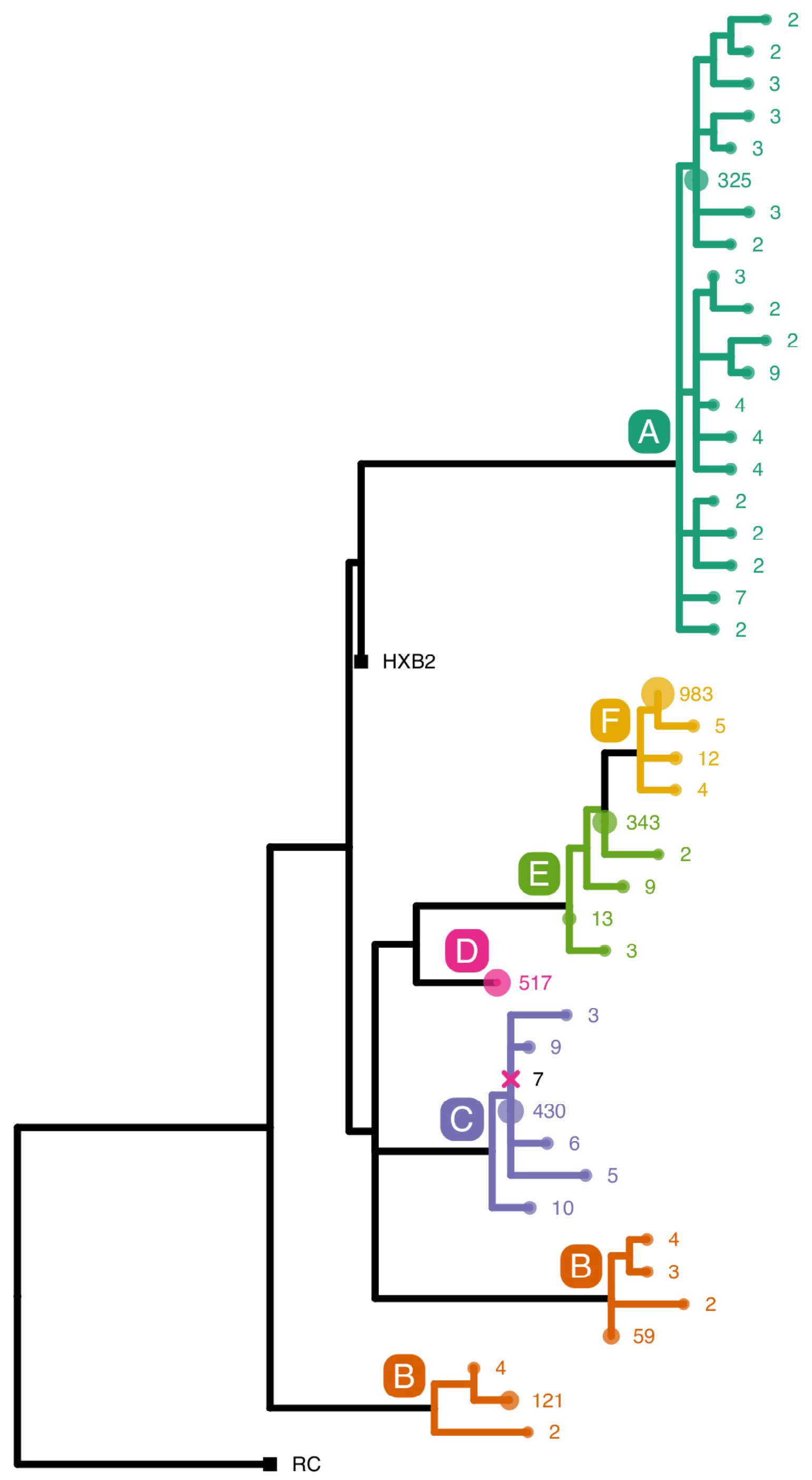
Wymant, Hall *et al.*, MBE 2017

Alignment with MAFFT: Katoh *et al.*, Nucleic Acids Res. 2002

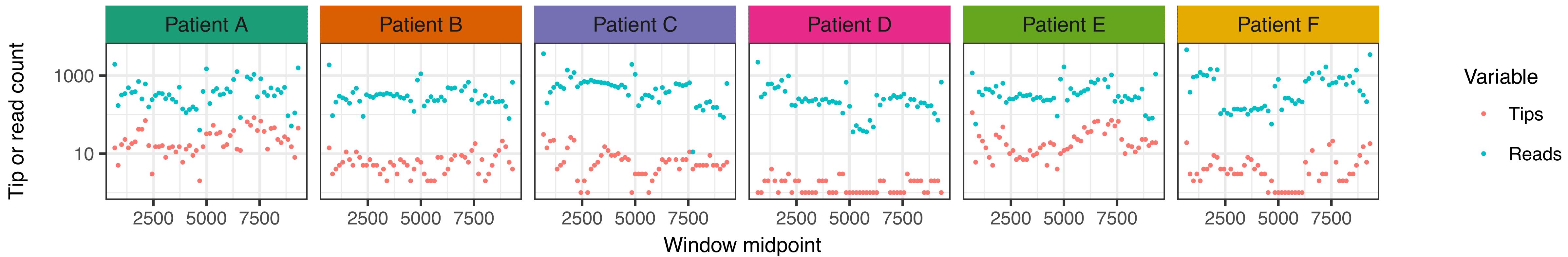
Phylogenetic inference with RAxML: Stamatakis, Bioinformatics 2014







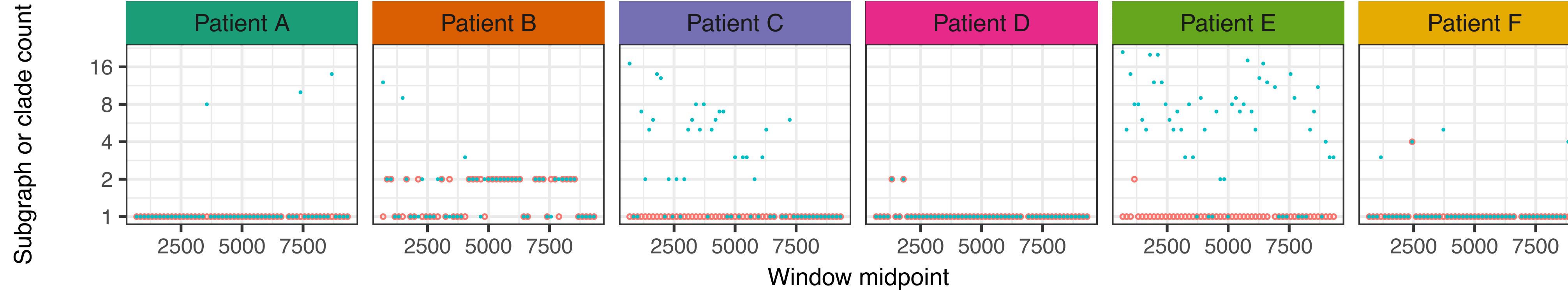
Per-patient summary statistics along the genome



100-1000 reads,
10-100 unique reads.
Diverse!

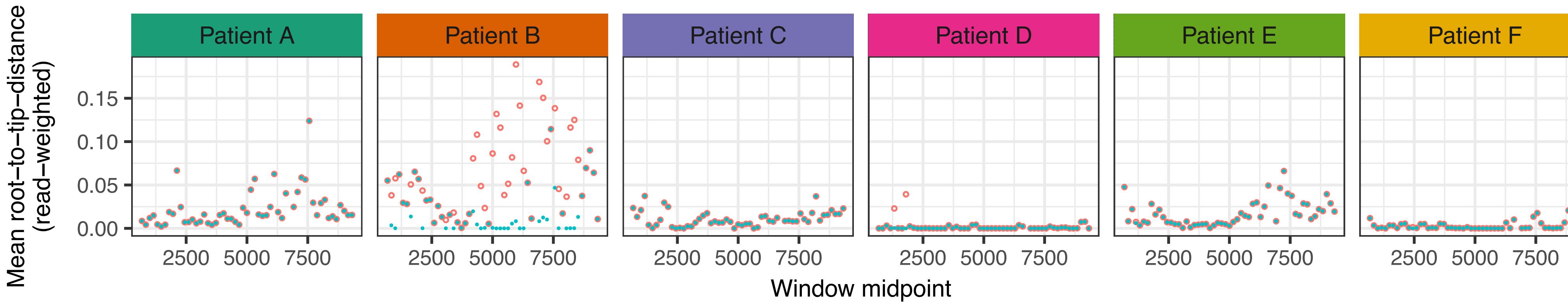


100-1000 reads,
1-2 unique reads.
Not diverse!



Variable

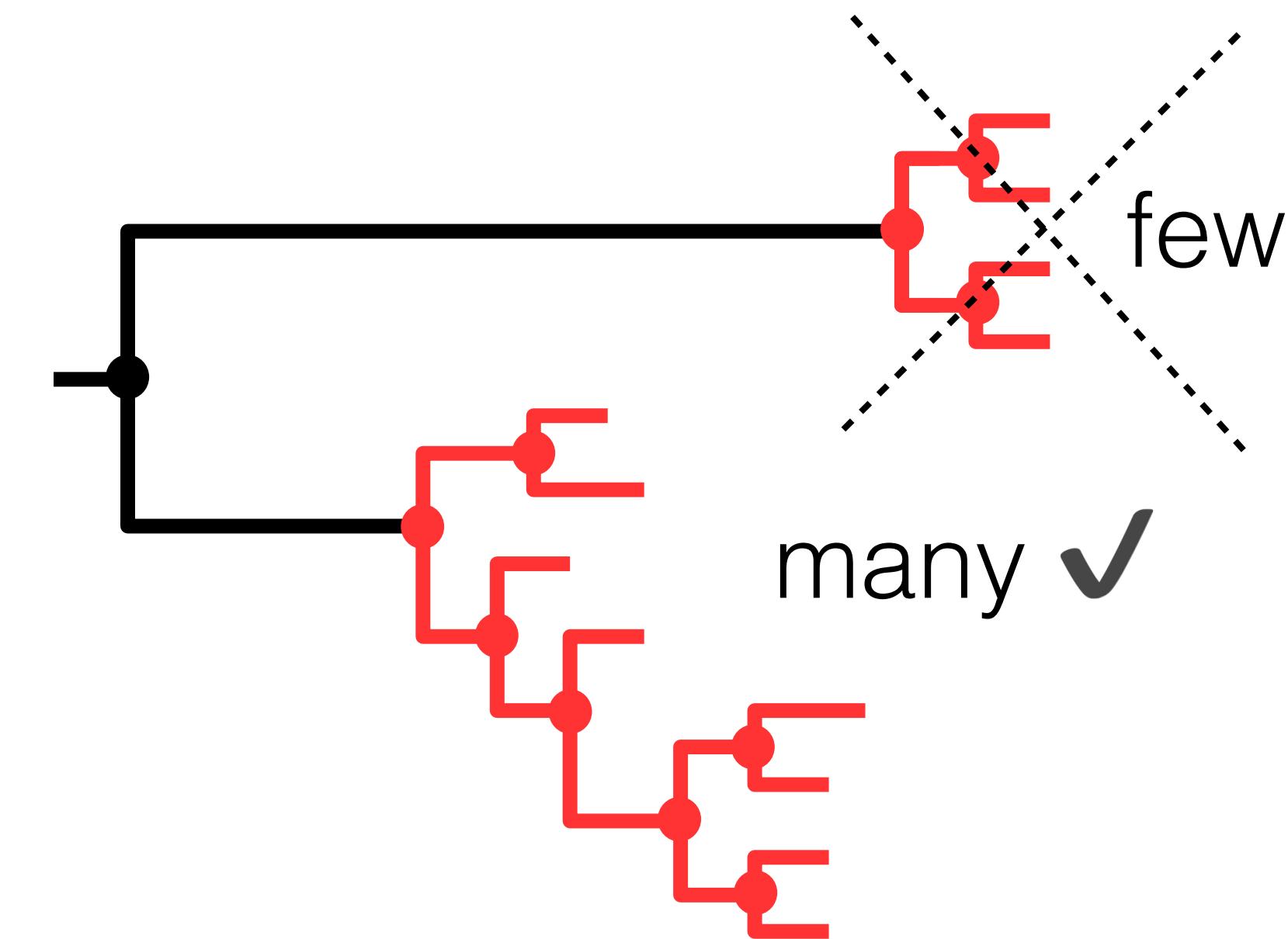
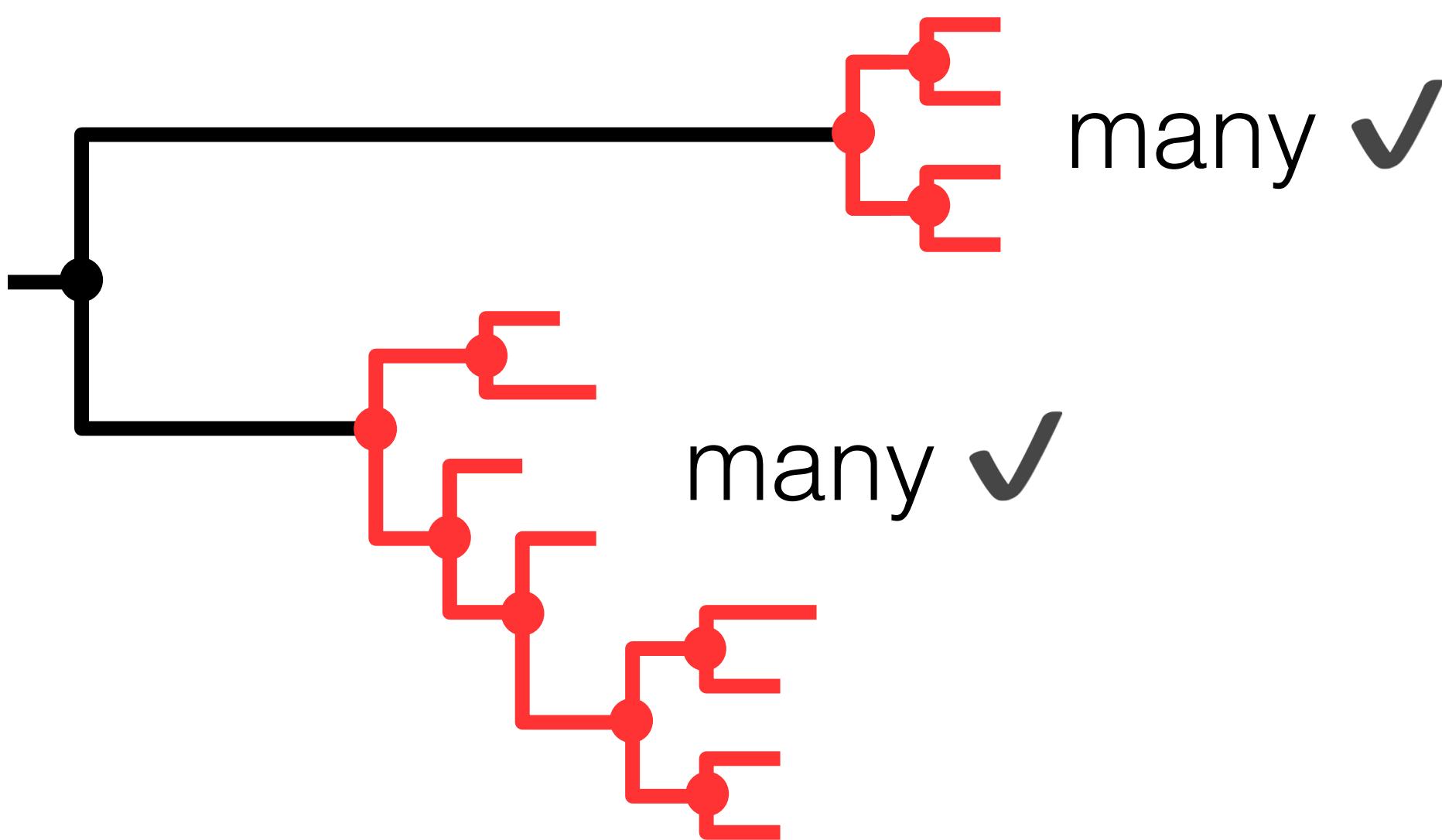
- Subgraphs
- Clades



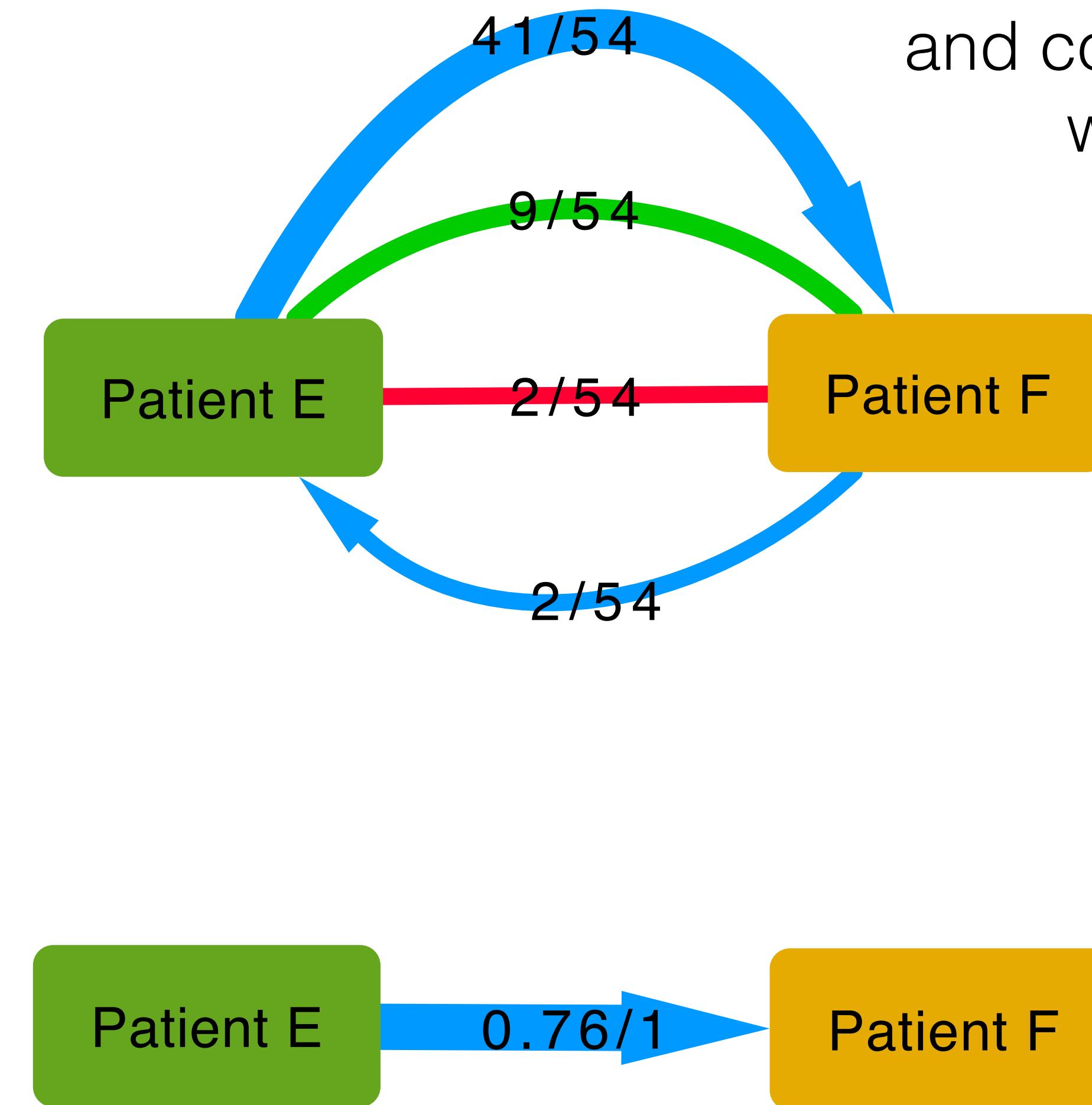
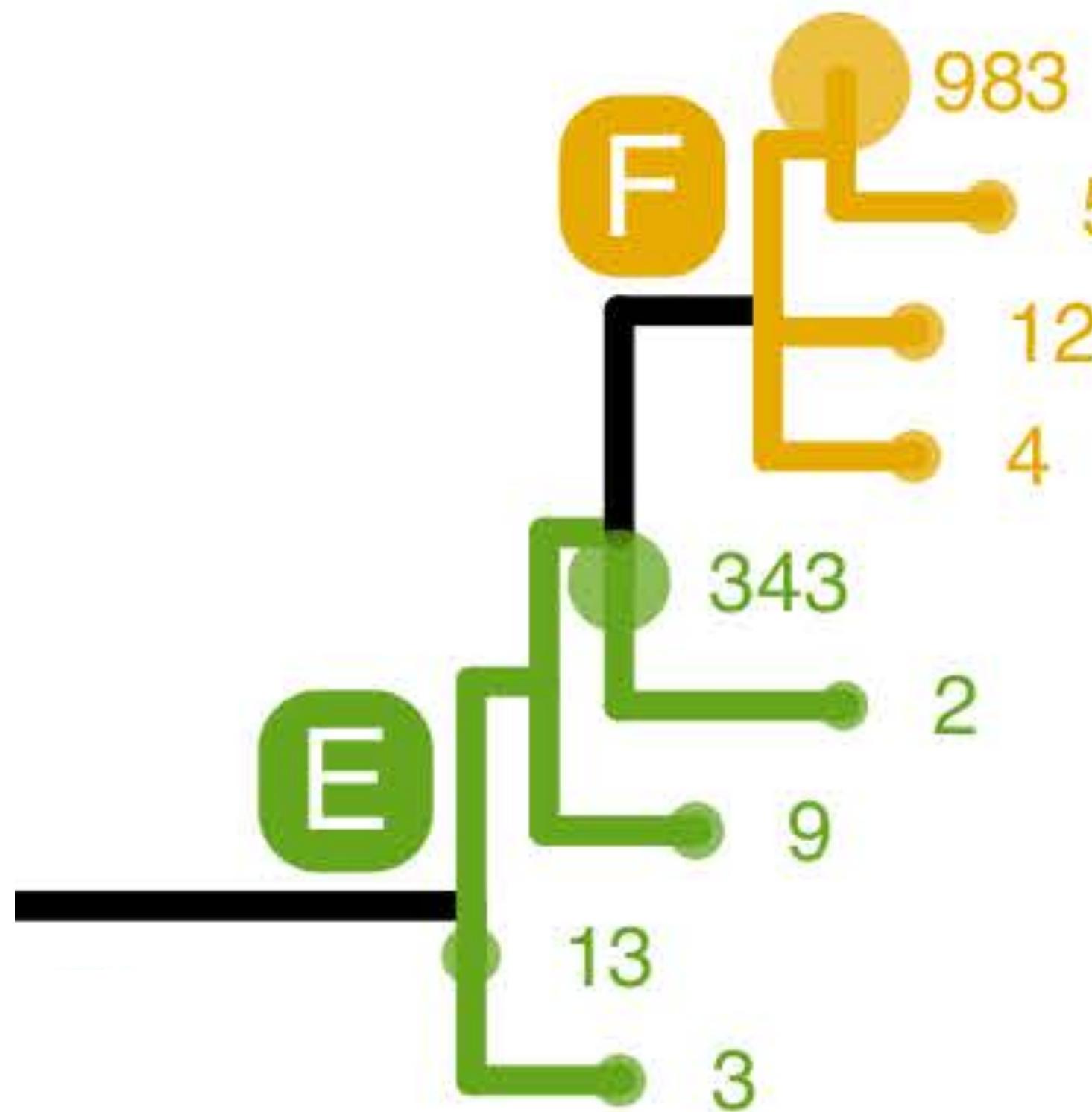
Tip set

- All
- Largest subgraph

Dually infected individual or contaminated sample?

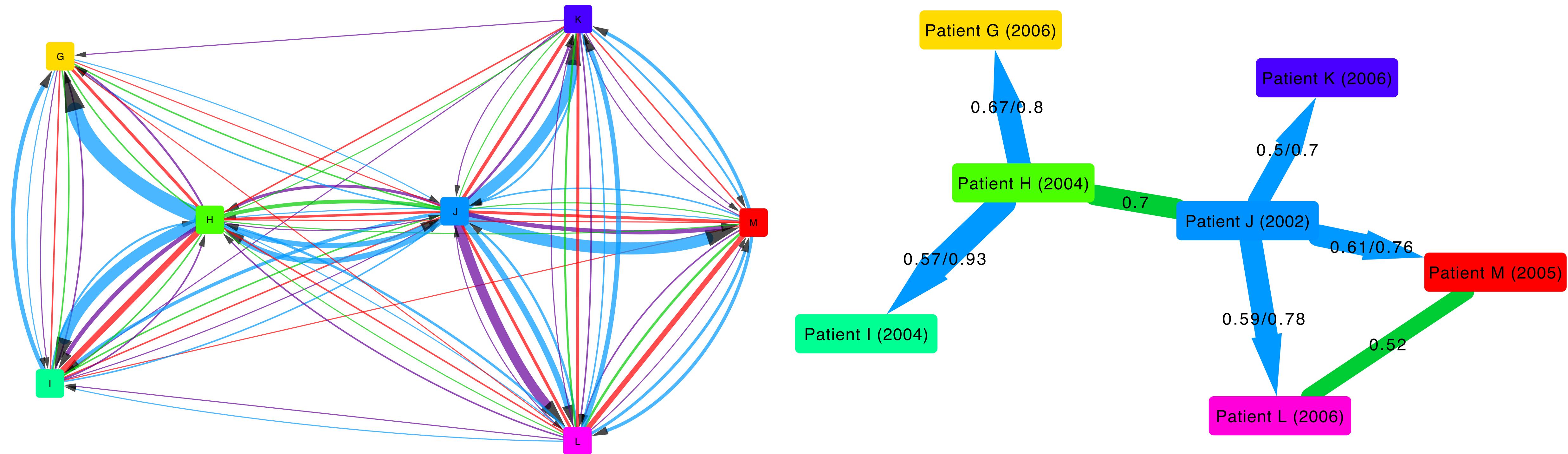


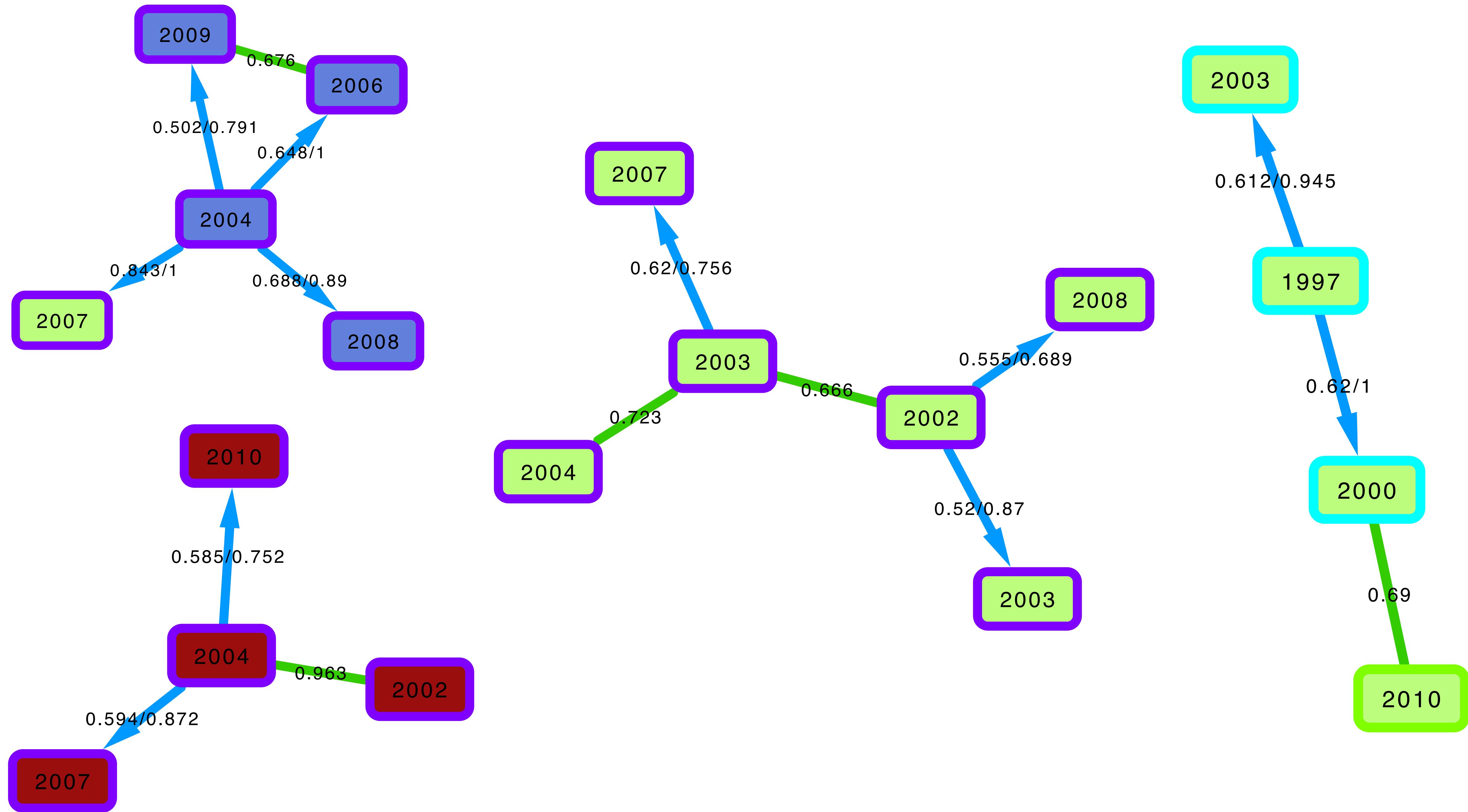
Transmission



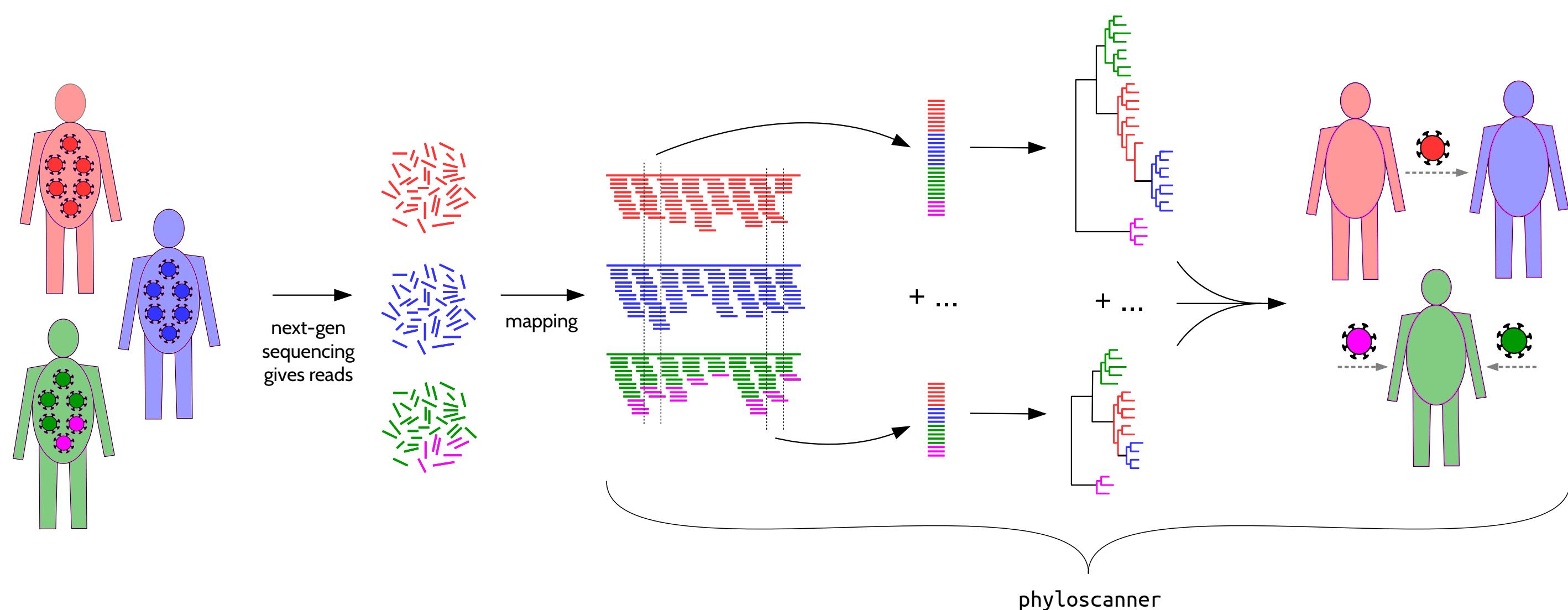
Summary: categorise
and count genomic
windows

Simpler summary: aggregate categories





Summary



phyloscanner creates within- & between-host phylogenies along the genome from NGS reads. Analysing these, or other such phylogenies, it identifies

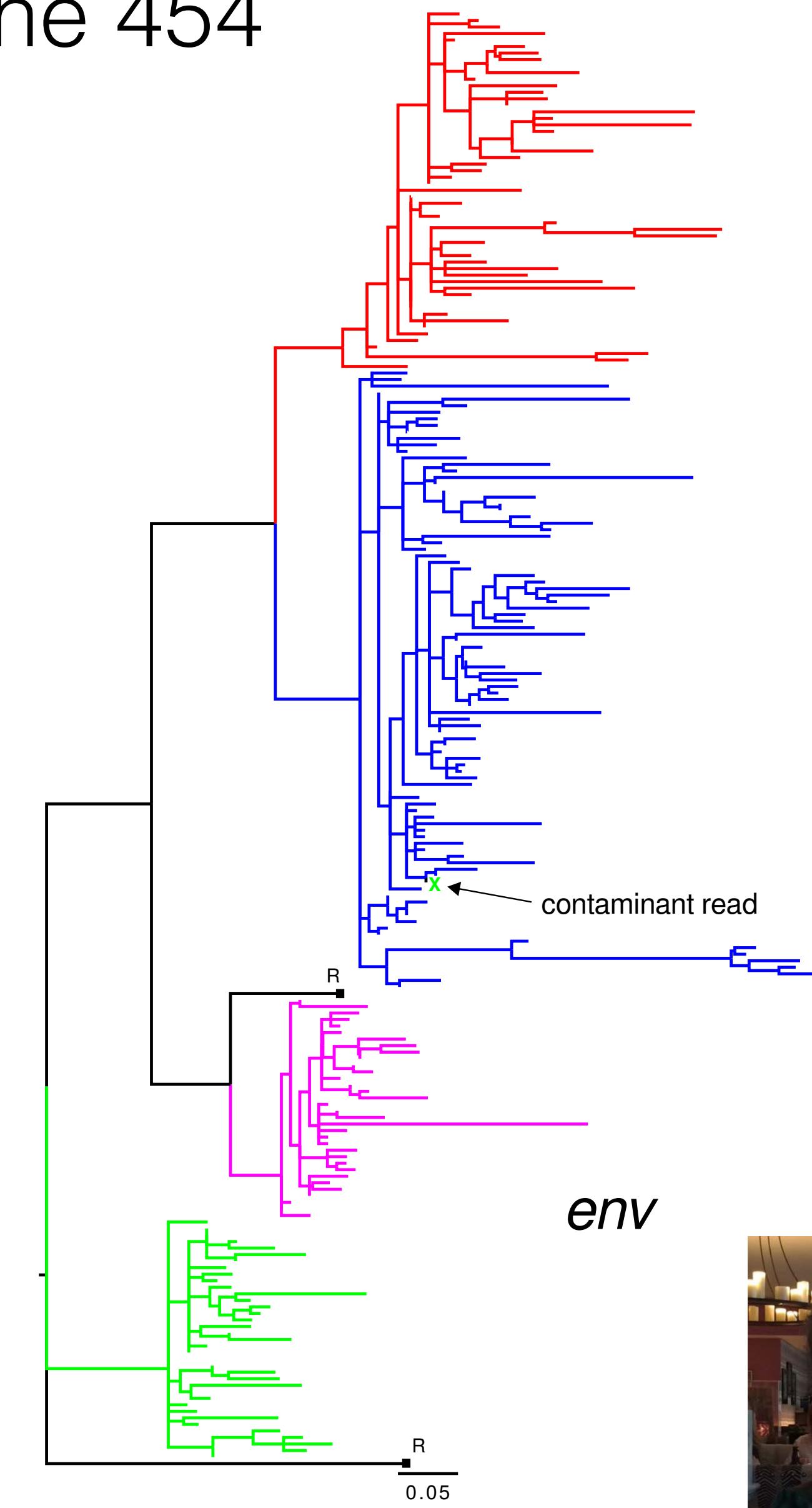
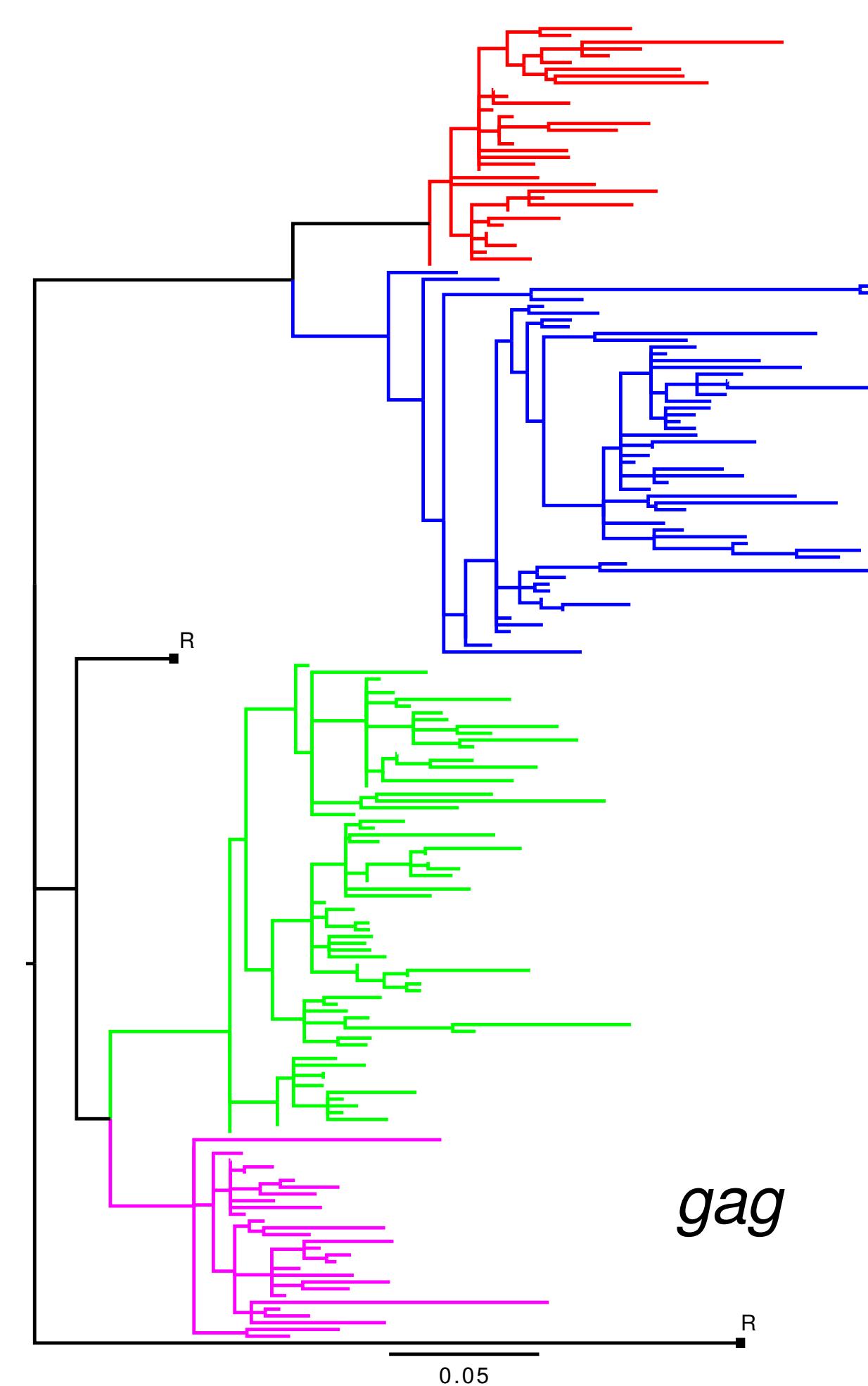
- **transmission**: one individual's pathogen population being ancestral to someone else's
- **multiple infection**: distinct subpopulations in the same host
- **contamination**: exact duplicates, phylogenetic outliers
- **recombination**

Wymant, Hall *et al.* MBE 2017

GitHub.com/BDI-pathogens/phyloscanner

Extra slides: application to other sequencing platforms and pathogens, and using phyloscanner

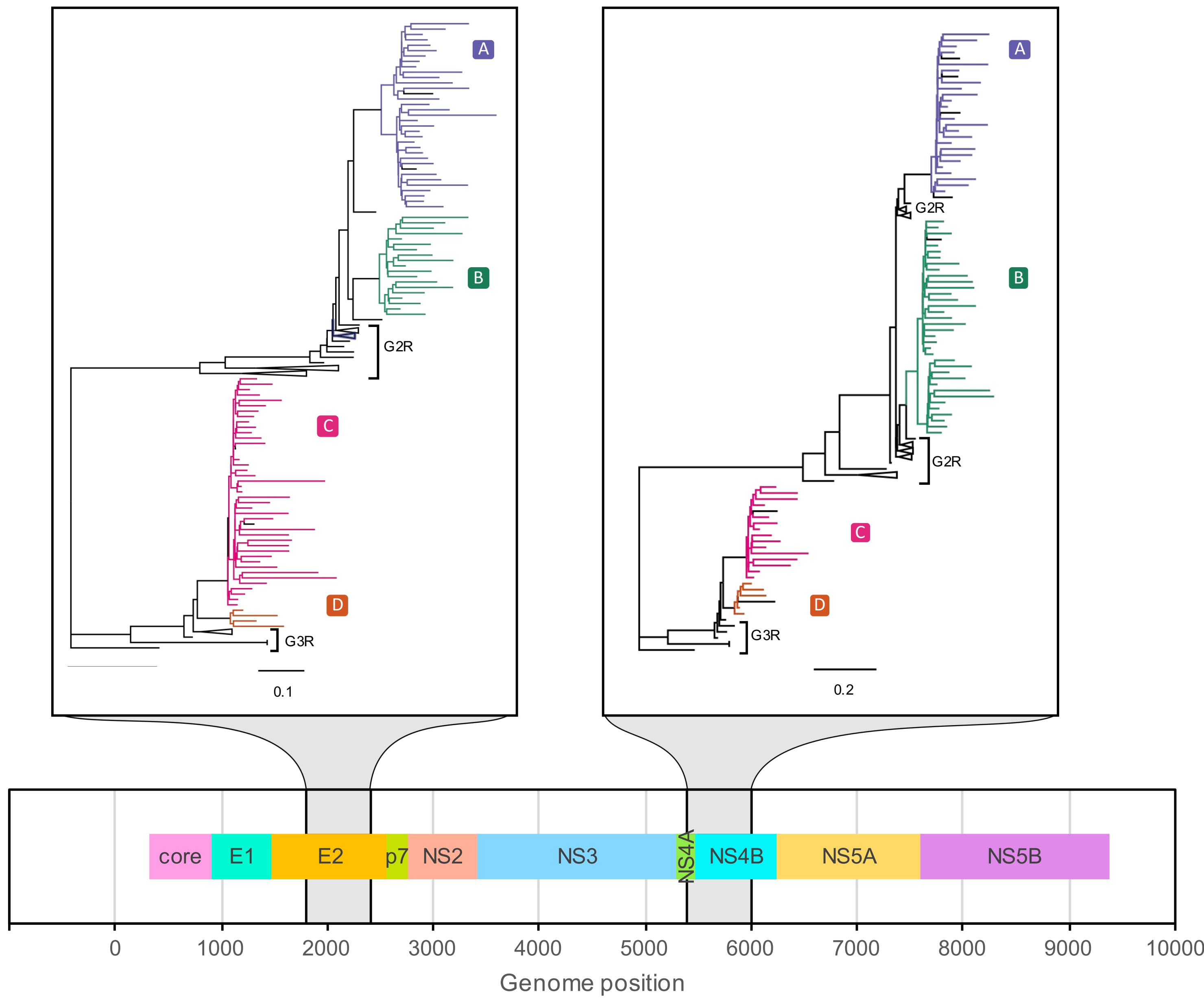
Four more BEEHIVE patients sequenced with Roche 454



Tanya Golubchik



Hepatitis C Virus Sequenced With Oxford Nanopore



The Stop-HCV Consortium:

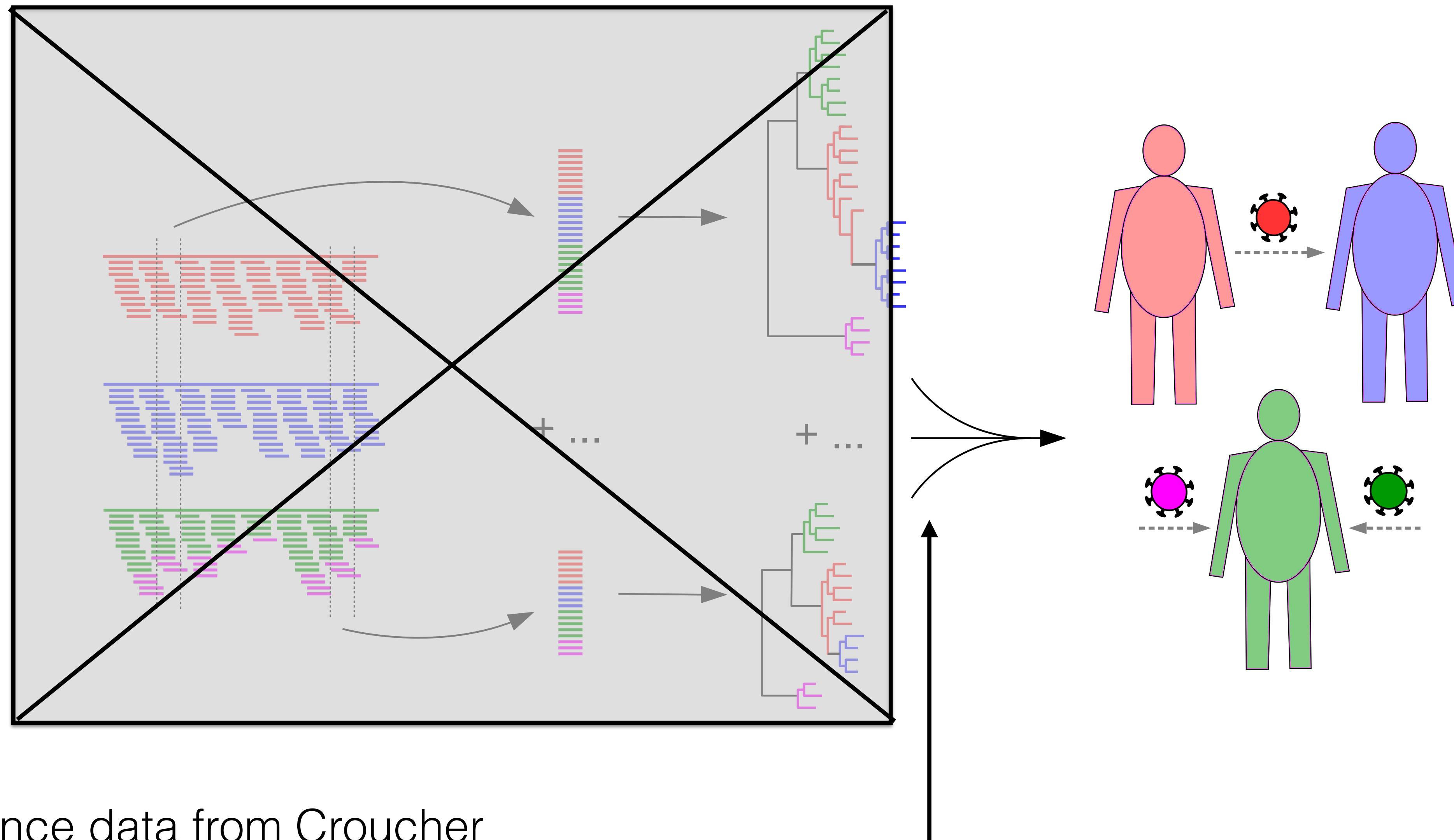
Eleanor Barnes, Diana Koletzki
Jonathan Ball Natasha Martin
Diana Brainard Benedetta Massetto
Gary Burgess Tamyo Mbisa
Graham Cooke John McHutchison
John Dillon Jane McKeating
Graham R Foster John McLauchlan
Charles Gore Alec Miners
Neil Guha Andrea Murray
Rachel Halford Peter Shaw
Cham Herath Peter Simmonds
Chris Holmes Chris C A Spencer
Anita Howe Paul Targett-Adams
Emma Hudson Emma Thomson
William Irving Peter Vickerman
Salim Khakoo Nicole Zitzmann
Paul Klenerman



David Bonsall



Mariateresa de Cesare



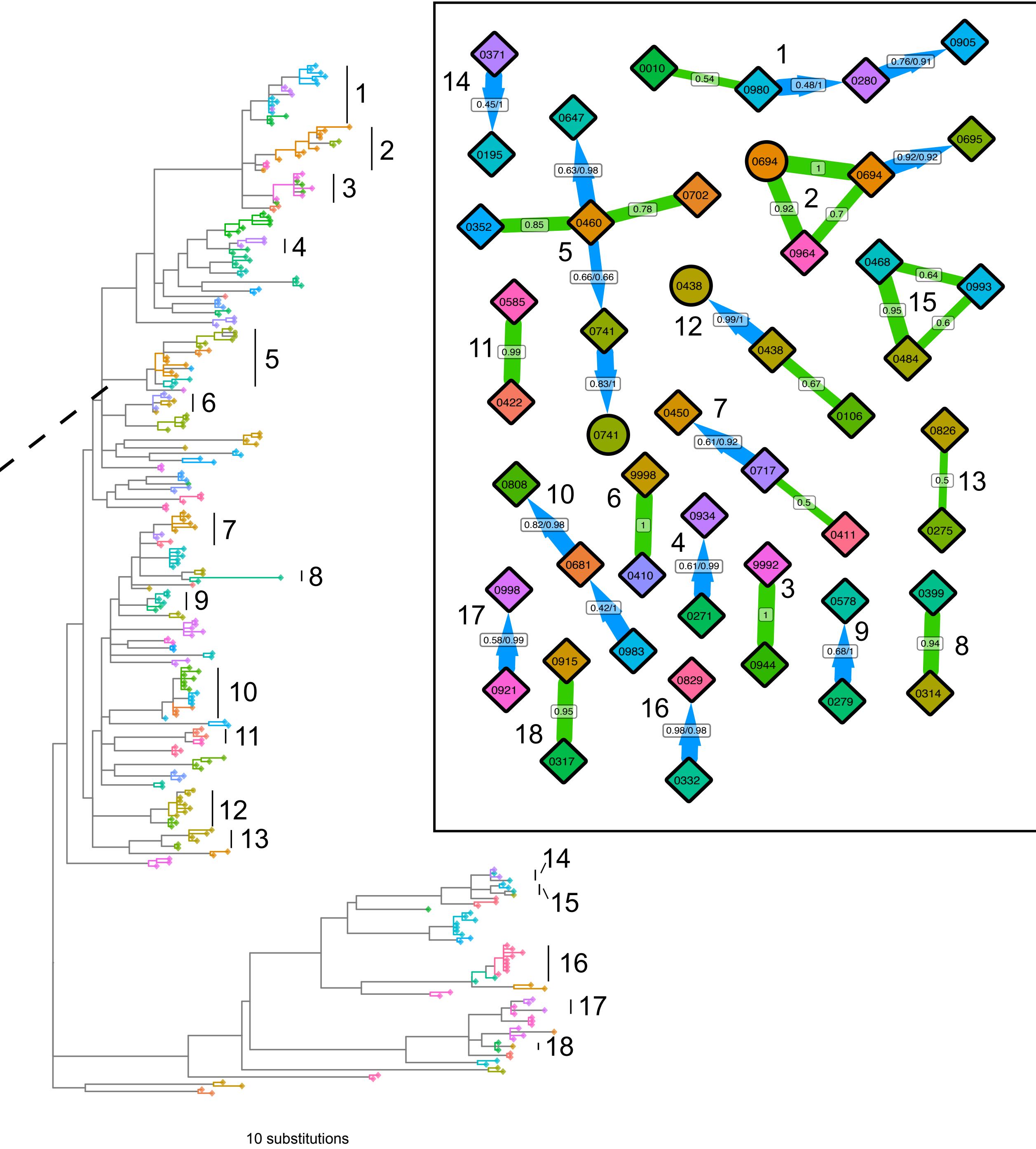
Sequence data from Croucher
et al. PLoS Bio. 2016: multiple
colony picks per carrier of
S. pneumoniae



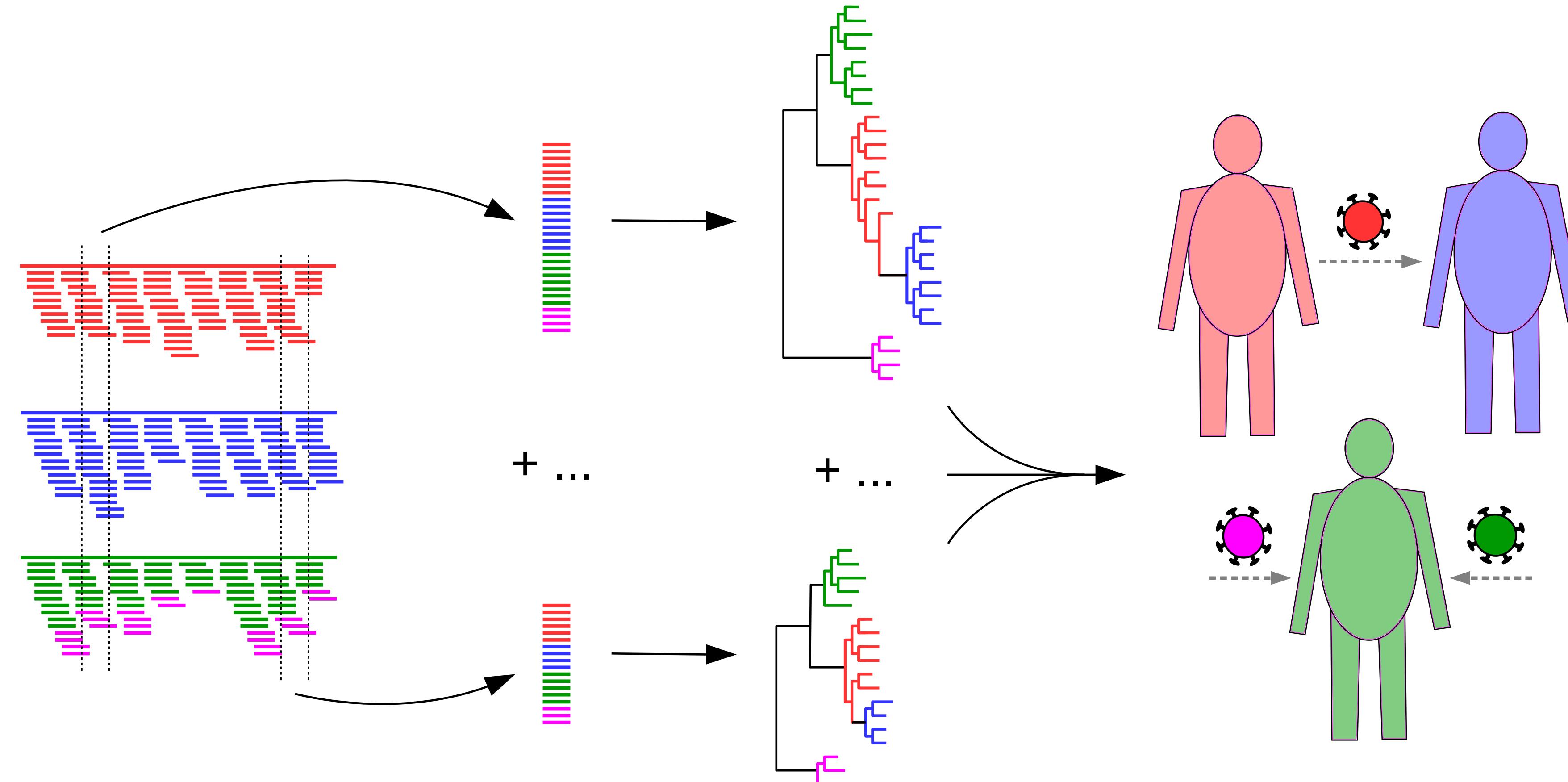
100 posterior trees
with MrBayes

The Maela Pneumococcal Collaboration:

Stephen Bentley
Claire Chewapreecha
Nicholas J. Croucher
Simon Harris
Jukka Corander
David Goldblatt
Julian Parkhill
Francois Nosten
Claudia Turner
Paul Turner



Using phyloscanner: [GitHub.com/BDI-pathogens/phyloscanner](https://github.com/BDI-pathogens/phyloscanner)

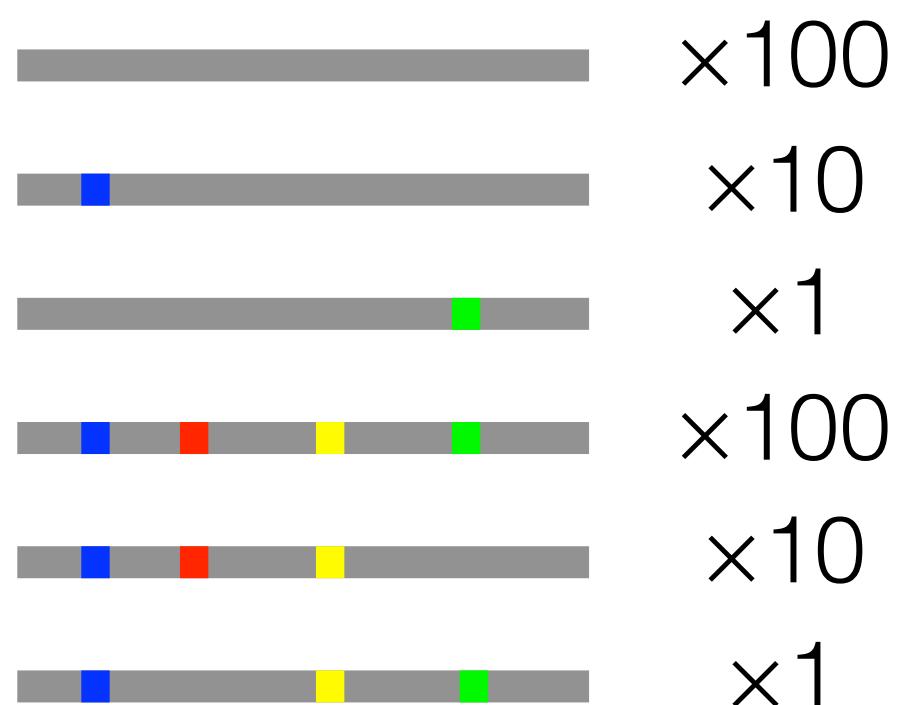
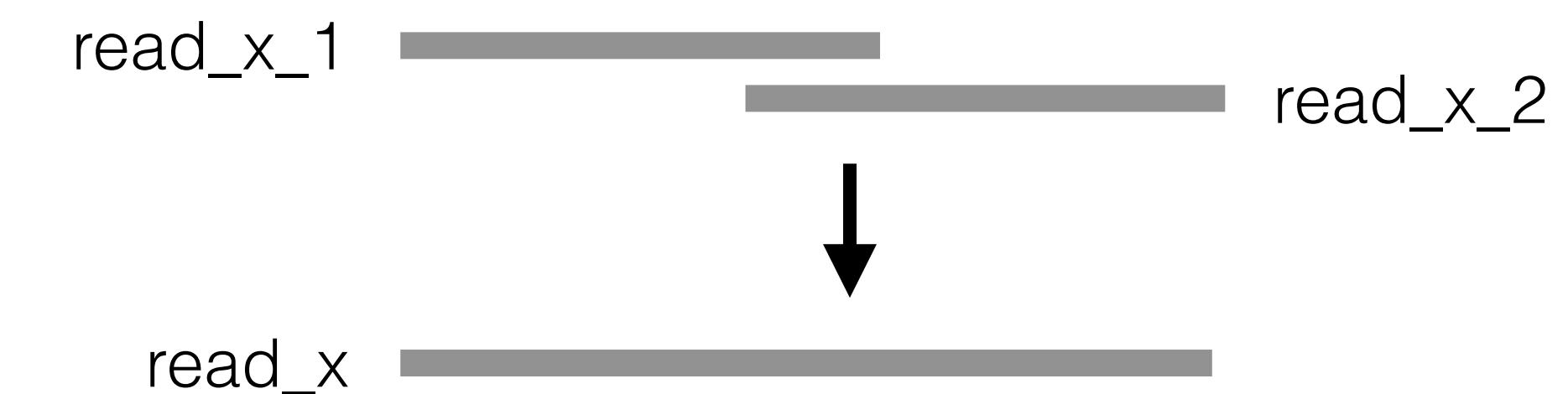


```
$ phyloscanner_make_trees.py ListOfBamFiles.csv --windows 1,300,301,600,...
```

```
$ phyloscanner_analyse_trees.R TreeFiles OutputString ChoiceOfHostStateReconstruction
```

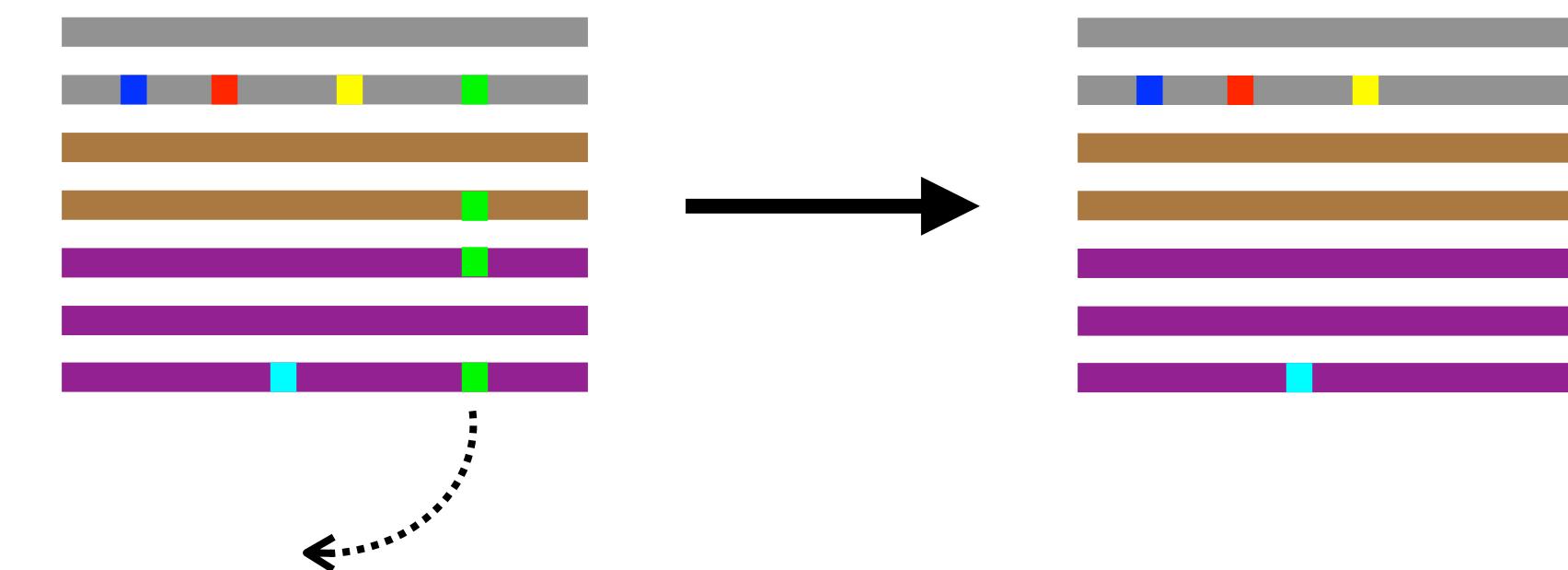
Options

Merge overlapping paired reads into longer reads.



Similarity- and frequency-based read merging, for speed.

Excise specified positions
(e.g. sites under selection).



Options

- Include known references with the reads
- Trim and/or discard low-quality reads
- Minimum read count
- Pass any RAxML options, e.g. bootstraps, model specs
- Choice of ancestral host-state reconstruction algorithms & parameters
- Transmission inference parameters, e.g. distance thresholds, distance normalisation over the genome

Trivially parallelisable: split the whole genome into windows, run each window as a separate job on your cluster (in addition to multithreading RAxML).