

Introduction to phylogenetics

Alexandra Blenkinsop

With thanks to Chris Wymant and Thibaut Jombart

February 24th 2024

- Why phylogenetics?
- How to interpret a phylogeny
- How to run a phylogenetic pipeline
 - Data, steps and tools for inferring ancestral relationships
- Analysing phylogenetic trees

What can phylogenetics tell us?

Phylogenetics for understanding epidemics

How can we target our interventions in an epidemic to have a greater impact?

We can learn about the epidemic through analysis of pathogen sequence data.

Examples:

- characterising spatial transmission dynamics
- characterising population-level drug resistance
- detection of new variants
- estimating epidemiological parameters
- identifying population-level drivers of transmission

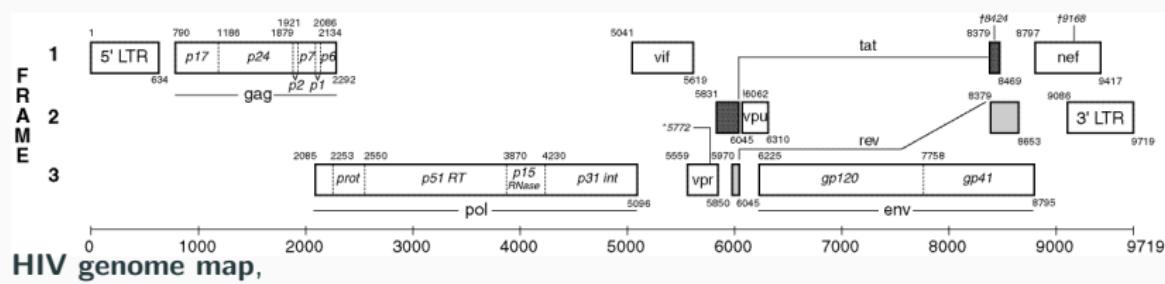
Viral sequences

Nucleotides are the building blocks of nucleic acids - chains of these (sequences) encode information in DNA/RNA

- Consensus sequence
 - One representative sequence of most frequent nucleotide at each position
- Next generation sequencing (NGS)
 - Sequence thousands of reads simultaneously
 - Can be used for inferring direction of transmission
 - Supports robustness of conclusions

Viral sequences

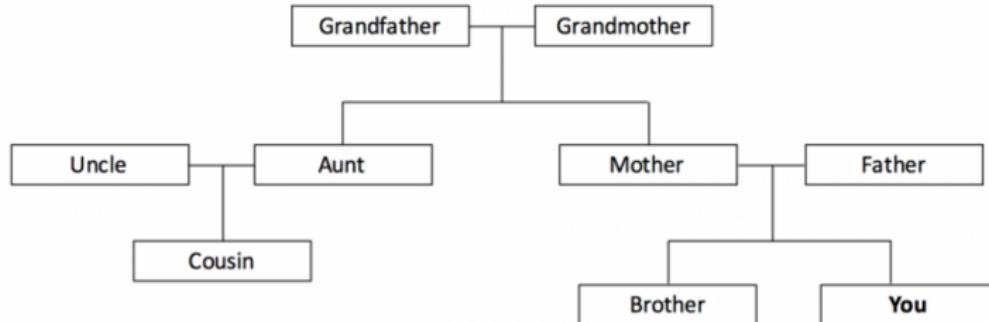
- Whole genome sequence covers all nucleotide positions
- Partial genome sequence covers part of the genome (e.g. *gag*, *pol*, *env* for HIV)



HIV genome map,

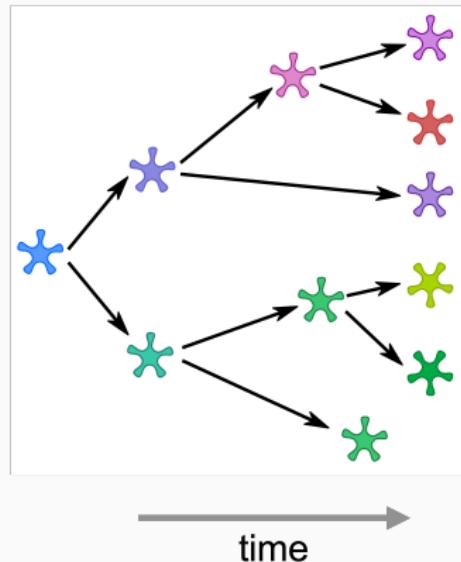
<https://www.hiv.lanl.gov/content/sequence/HIV/MAP/landmark.html>

Using sequence data to infer ancestral relationships

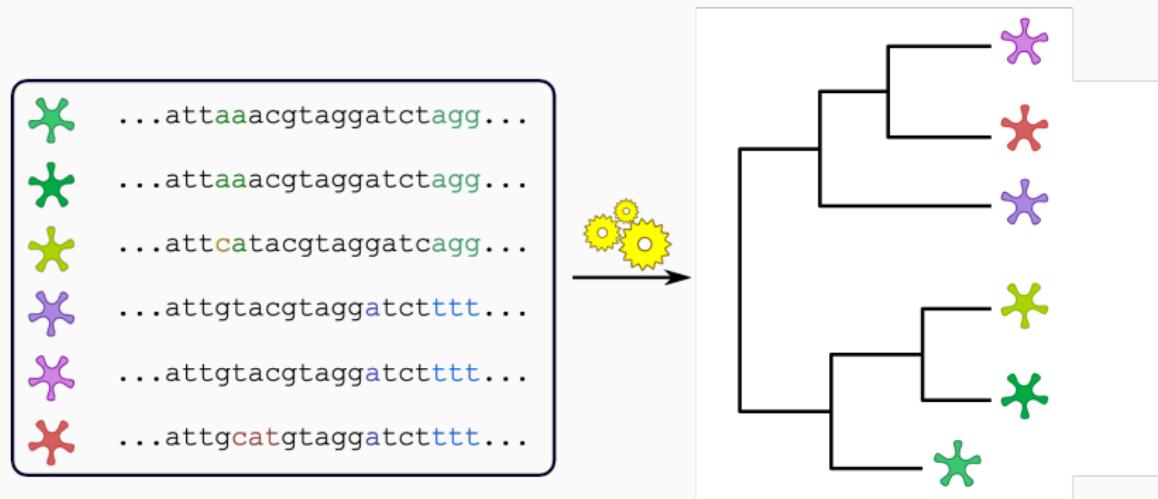


Genetic changes (mutations) accumulate over time

- Base substitution: replacement of a nucleotide (e.g. a→t)
- Insertions: extra nucleotide added (e.g. at →agt)
- Deletions: omission of a nucleotide during replication (e.g. agt →at)

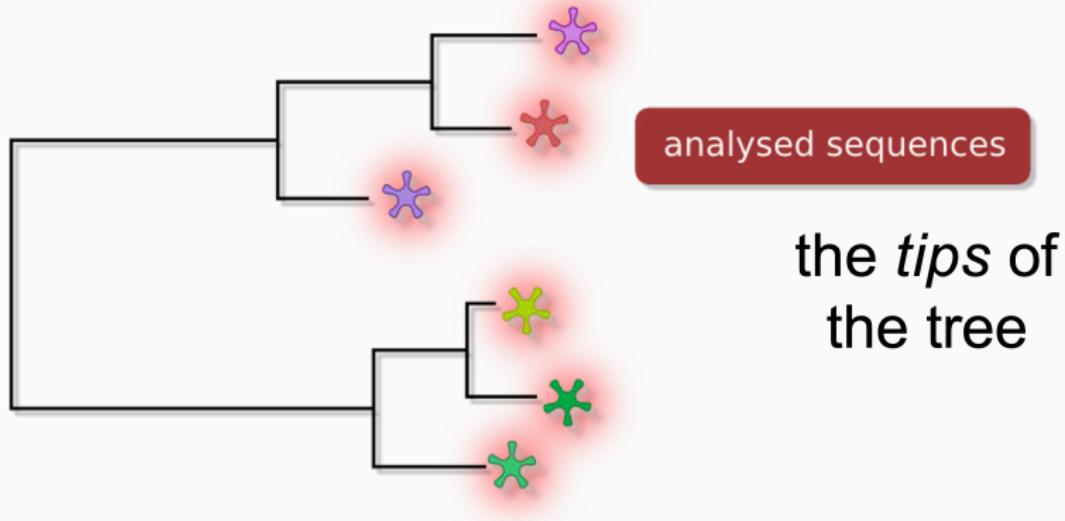


Using substitution patterns to reconstruct the evolutionary history



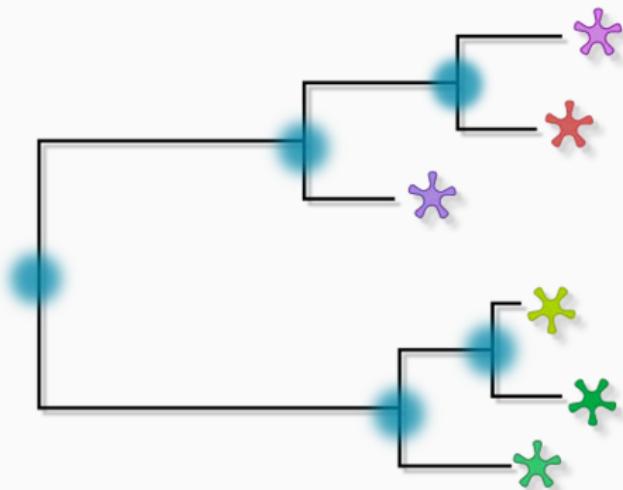
Interpreting a phylogeny

Features of a phylogeny



Most Recent Common Ancestors
(MRCA)

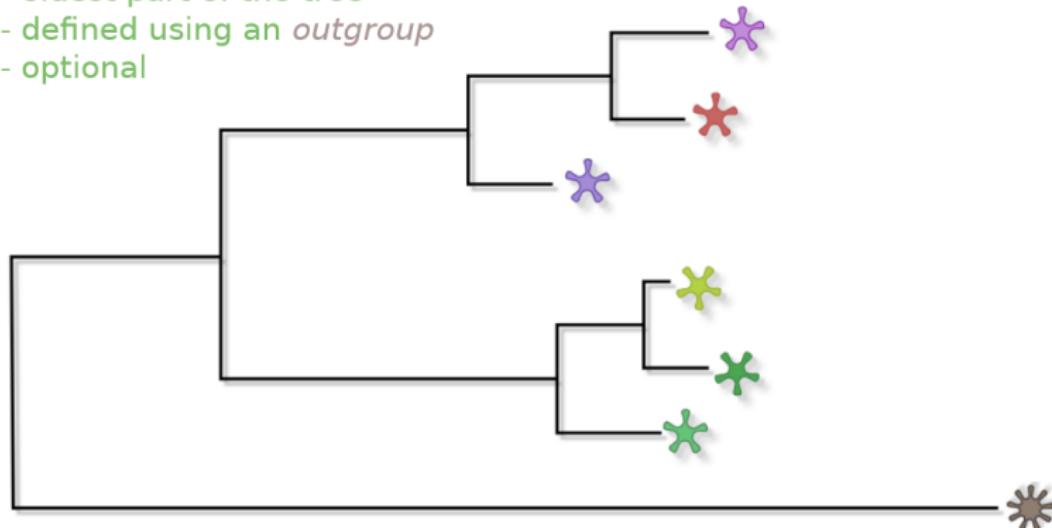
*the nodes
of the tree*



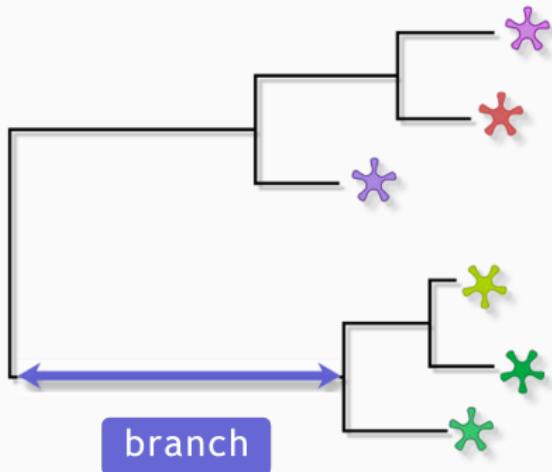
Features of a phylogeny

Root

- oldest part of the tree
- defined using an *outgroup*
- optional



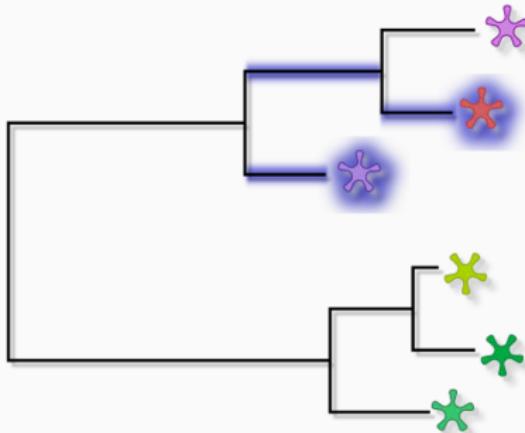
Features of a phylogeny



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

Features of a phylogeny

Small distance \Rightarrow small amount of evolution, from which we infer epidemiological proximity of pathogens



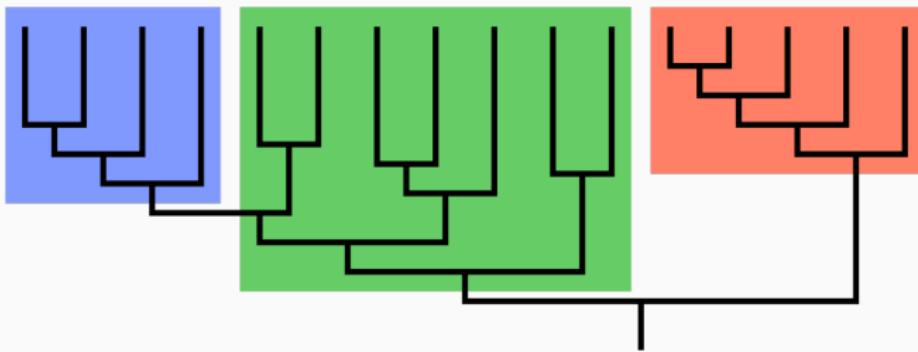
distances between tips

- "patristic" distance: sum of branch lengths

Features of a phylogeny

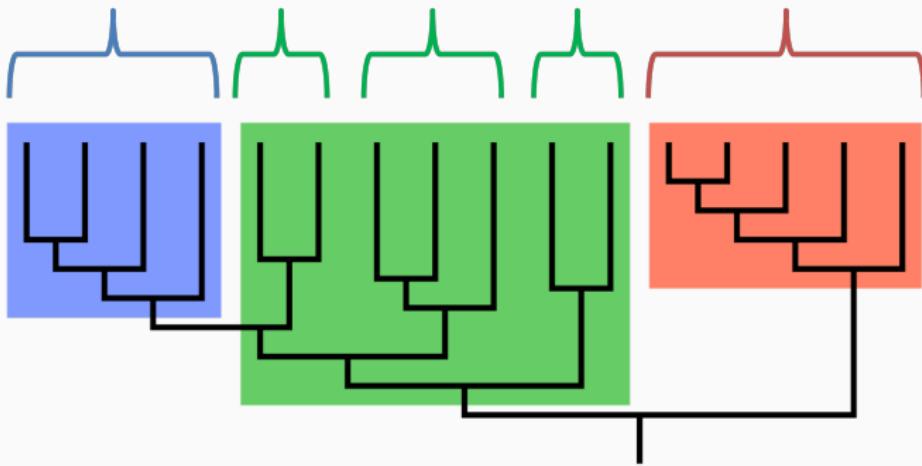
Clade: A common ancestor and all its descendants (monophyletic group)

Which of these groups are clades?



Features of a phylogeny

Clade: A common ancestor and all its descendants (monophyletic group)



Phylogenetic algorithms infer the phylogeny using an assumed substitution model

Main approaches:

- Distance-based e.g. neighbour-joining
- Maximum parsimony
- Maximum-likelihood
- Bayesian

Substitution models for phylogenetic inference

Parametric approaches are based on molecular models of how the pathogen accumulates substitutions (e.g. the relative rate of A→C versus T→G, etc.)

Non-exhaustive examples:

- Jukes and Cantor
 - equal base frequencies
 - equal rates of mutation
 - one parameter (substitution rate)
- Kimura and Nei
 - relaxes second assumption to reflect difference in rate of transition and transversion mutations
 - two parameters for rate
- Hasegawa–Kishino–Yano (HKY)
 - nucleotides occur at different frequencies
 - transitions and transversions occur at different rates
- General Time Reversible (GTR)
 - nucleotides occur at different frequencies
 - different rates of substitution for each pair of nucleotides

Interpreting a phylogeny

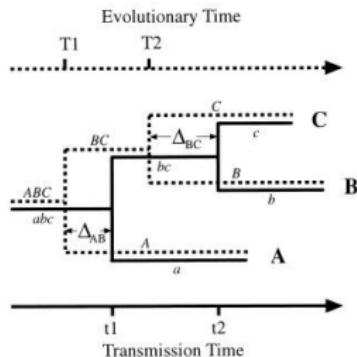
Tree summary statistics can help characterise topological features and ancestral relationships, e.g.

- Max tree depth
- Maximum tree width
- Phylogenetic diversity
- Node to tip (bifurcation) ratio

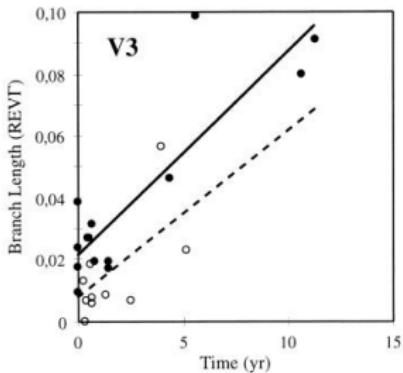
Interpreting a phylogeny

Genetic divergence (subst/site) =
evolutionary rate (substs/site/year) \times
divergence time (years)

A



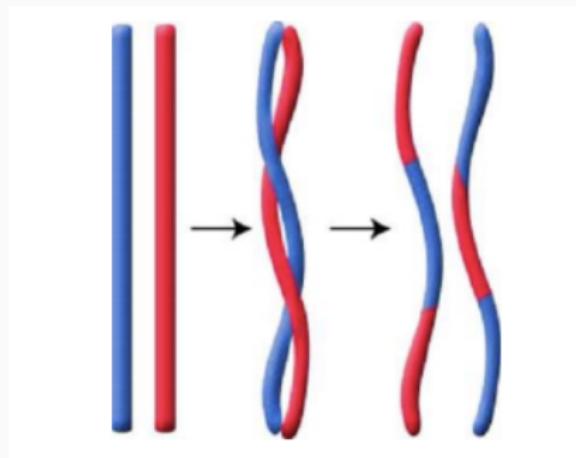
B



Processes which can affect phylogenetic inference

Recombination

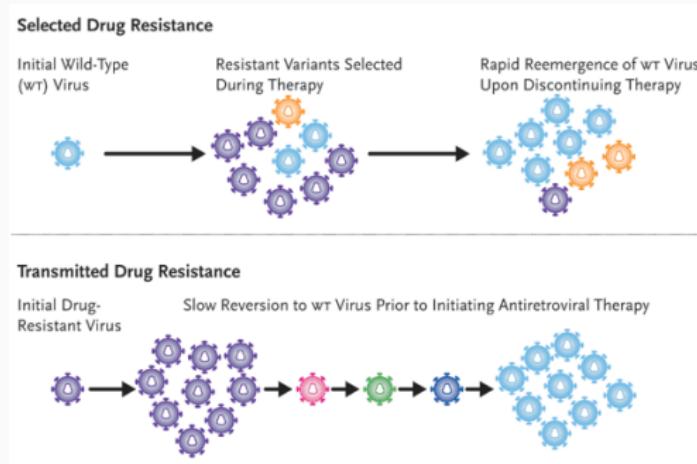
- When genome segments from different viruses are spliced together
- Genetic diversity no longer depends only on evolutionary rates and time



Processes which can affect phylogenetic inference

Drug resistance

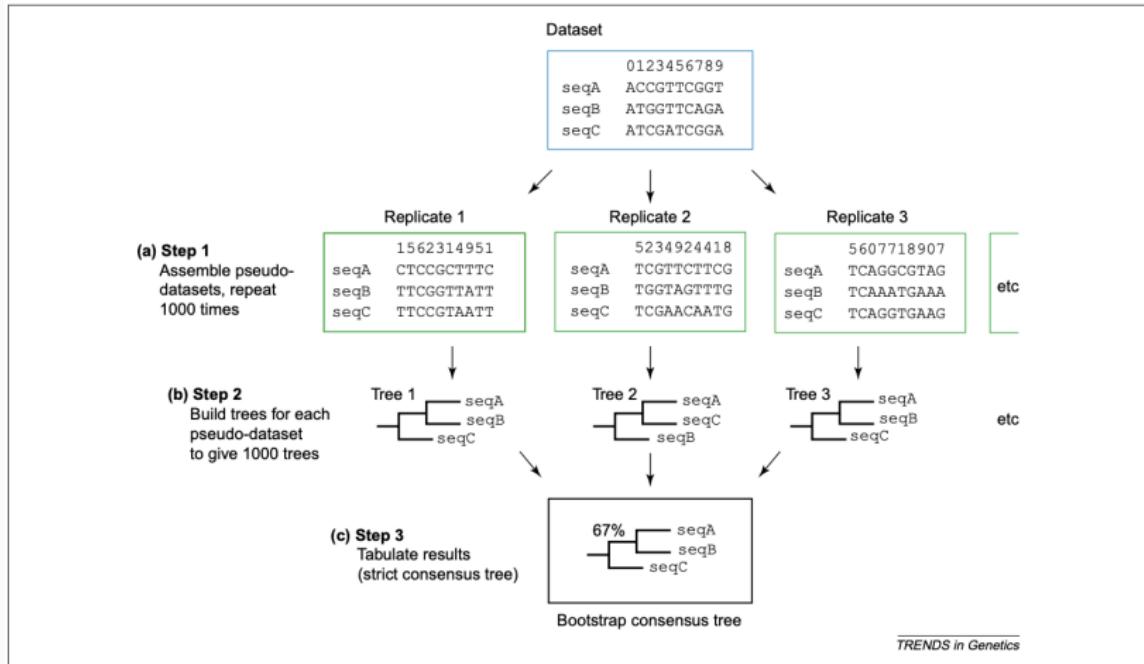
- Drug resistant mutations (DRMs) counteract treatment-mediated inhibition of viral replication
- Are the result of drug-selective pressure, either through treatment or can be transmitted
- Can bias phylogenetic inference



Kuritzkes, 2004

Quantifying uncertainty

Estimating uncertainty in inferred tree topology via bootstrapping



Baldauf, 2003 *Trends in Genetics*

Limitations of phylogenetic analyses

- Inferences are uncertain, e.g. in the case of pathogen spread, we are approximating the true unknown transmission tree
- Sampling bias

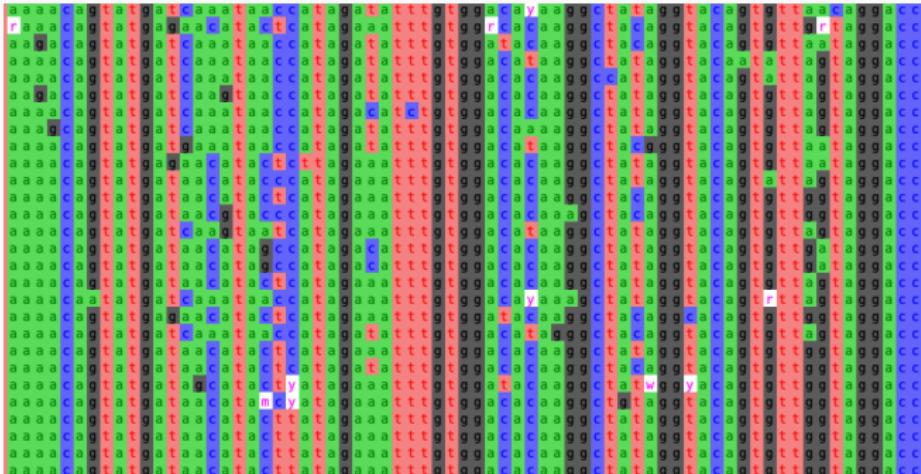
Statistical modelling can help by:

- Incorporating phylogenetic uncertainty
- Accounting for the unsampled population to make population-level inferences
- Amalgamating other data sources (e.g. clinical, mobility, contact data) to understand epidemiological trends

Running a phylogenetic pipeline

Sequence alignment

- Align sequences using an algorithm which seeks to minimise some distance measure
- Pairwise vs multiple sequence alignment
- Software: *MAFFT*, *CLUSTALW*, *virulign*



Mask drug resistant mutations (optional)

Mask known DRMs in the alignment to minimise bias in inferred ancestral relationships in phylogeny from common drug-resistant sites

A sequence alignment of DNA or RNA. The sequences are color-coded by base: A (blue), T (red), C (green), G (purple). Some positions are highlighted with black boxes. The first sequence has a black box around the 10th position. The second sequence has a black box around the 10th position. The third sequence has a black box around the 10th position. The fourth sequence has a black box around the 10th position. The fifth sequence has a black box around the 10th position. The sixth sequence has a black box around the 10th position. The seventh sequence has a black box around the 10th position. The eighth sequence has a black box around the 10th position. The ninth sequence has a black box around the 10th position. The tenth sequence has a black box around the 10th position. The eleventh sequence has a black box around the 10th position. The twelfth sequence has a black box around the 10th position. The thirteenth sequence has a black box around the 10th position. The fourteenth sequence has a black box around the 10th position. The fifteenth sequence has a black box around the 10th position. The sixteenth sequence has a black box around the 10th position. The seventeenth sequence has a black box around the 10th position. The eighteenth sequence has a black box around the 10th position. The nineteenth sequence has a black box around the 10th position. The twentieth sequence has a black box around the 10th position. The twenty-first sequence has a black box around the 10th position. The twenty-second sequence has a black box around the 10th position. The twenty-third sequence has a black box around the 10th position. The twenty-fourth sequence has a black box around the 10th position. The twenty-fifth sequence has a black box around the 10th position. The twenty-sixth sequence has a black box around the 10th position. The twenty-seventh sequence has a black box around the 10th position. The twenty-eighth sequence has a black box around the 10th position. The twenty-ninth sequence has a black box around the 10th position. The thirtieth sequence has a black box around the 10th position. The thirty-first sequence has a black box around the 10th position. The thirty-second sequence has a black box around the 10th position. The thirty-third sequence has a black box around the 10th position. The thirty-fourth sequence has a black box around the 10th position. The thirty-fifth sequence has a black box around the 10th position. The thirty-sixth sequence has a black box around the 10th position. The thirty-seventh sequence has a black box around the 10th position. The thirty-eighth sequence has a black box around the 10th position. The thirty-ninth sequence has a black box around the 10th position. The forty-second sequence has a black box around the 10th position. The forty-third sequence has a black box around the 10th position. The forty-fourth sequence has a black box around the 10th position. The forty-fifth sequence has a black box around the 10th position. The forty-sixth sequence has a black box around the 10th position. The forty-seventh sequence has a black box around the 10th position. The forty-eighth sequence has a black box around the 10th position. The forty-ninth sequence has a black box around the 10th position. The fifty-second sequence has a black box around the 10th position. The fifty-third sequence has a black box around the 10th position. The fifty-fourth sequence has a black box around the 10th position. The fifty-fifth sequence has a black box around the 10th position. The fifty-sixth sequence has a black box around the 10th position. The fifty-seventh sequence has a black box around the 10th position. The fifty-eighth sequence has a black box around the 10th position. The fifty-ninth sequence has a black box around the 10th position. The sixty-second sequence has a black box around the 10th position. The sixty-third sequence has a black box around the 10th position. The sixty-fourth sequence has a black box around the 10th position. The sixty-fifth sequence has a black box around the 10th position. The sixty-sixth sequence has a black box around the 10th position. The sixty-seventh sequence has a black box around the 10th position. The sixty-eighth sequence has a black box around the 10th position. The sixty-ninth sequence has a black box around the 10th position. The seventy-second sequence has a black box around the 10th position. The seventy-third sequence has a black box around the 10th position. The seventy-fourth sequence has a black box around the 10th position. The seventy-fifth sequence has a black box around the 10th position. The seventy-sixth sequence has a black box around the 10th position. The seventy-seventh sequence has a black box around the 10th position. The seventy-eighth sequence has a black box around the 10th position. The seventy-ninth sequence has a black box around the 10th position. The eighty-second sequence has a black box around the 10th position. The eighty-third sequence has a black box around the 10th position. The eighty-fourth sequence has a black box around the 10th position. The eighty-fifth sequence has a black box around the 10th position. The eighty-sixth sequence has a black box around the 10th position. The eighty-seventh sequence has a black box around the 10th position. The eighty-eighth sequence has a black box around the 10th position. The eighty-ninth sequence has a black box around the 10th position. The ninety-second sequence has a black box around the 10th position. The ninety-third sequence has a black box around the 10th position. The ninety-fourth sequence has a black box around the 10th position. The ninety-fifth sequence has a black box around the 10th position. The ninety-sixth sequence has a black box around the 10th position. The ninety-seventh sequence has a black box around the 10th position. The ninety-eighth sequence has a black box around the 10th position. The ninety-ninth sequence has a black box around the 10th position. The one-hundredth sequence has a black box around the 10th position.

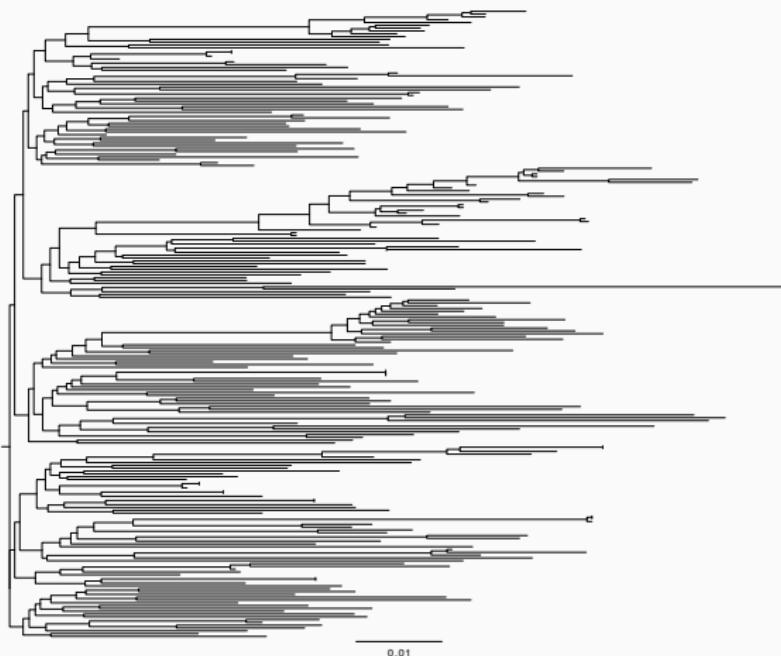
Build trees

Software: *FastTree*, *RaxML*, *IQTree*, *PhyML*, *MrBayes*, *BEAST*, *ape* (*R* package)



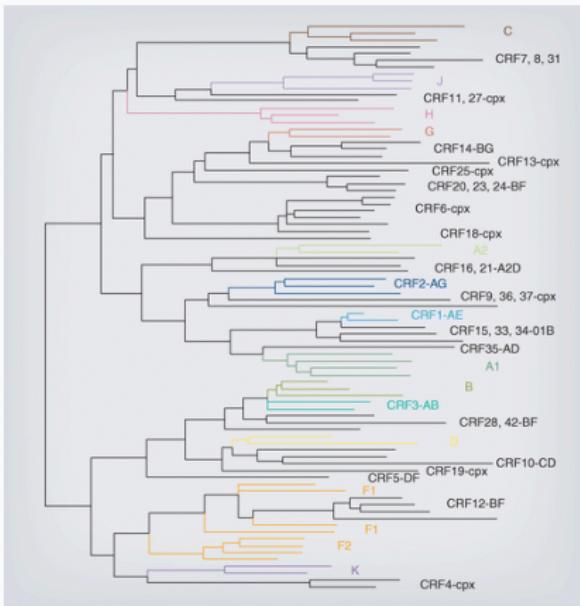
Check trees (if feasible)

- Visually inspect trees
- Very long branches may indicate a problem with the alignment



Select outgroups and root tree

Typically use reference sequences from a sequence database (e.g. GenBank)



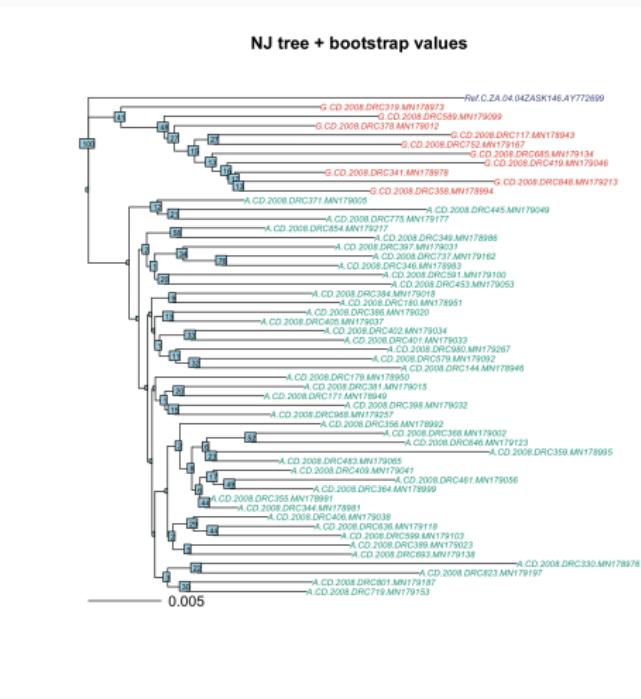
Castro-Nallar et al. (2012). *Future Virology*

Quantify uncertainty

- Bootstrap resample alignment with replacement
 - Build trees
 - Label internal nodes of central alignment with bootstrap values

Tools:

- boot.phylo in ape package in R
 - infer bootstrap trees in parallel using high-performance computing



Analysing phylogenetic trees

- Cluster analysis → how many, how big, importations
- Ancestral state reconstruction
 - Characterising epidemiological transmission dynamics → source attribution
 - Phylogeography → spatial transmission dynamics
- Phylodynamics → estimating population-level parameters which shape phylogenies

Identifying clusters

Distance-based clustering groups taxa with strong evidence of small genetic diversity, through user-specified bootstrap support and patristic distance thresholds

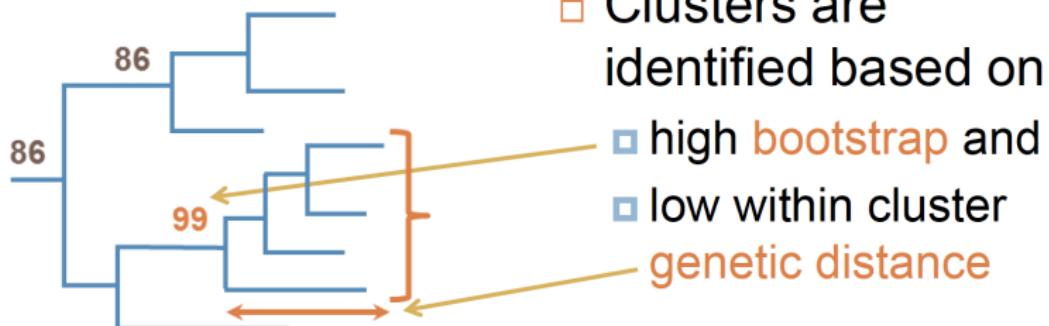
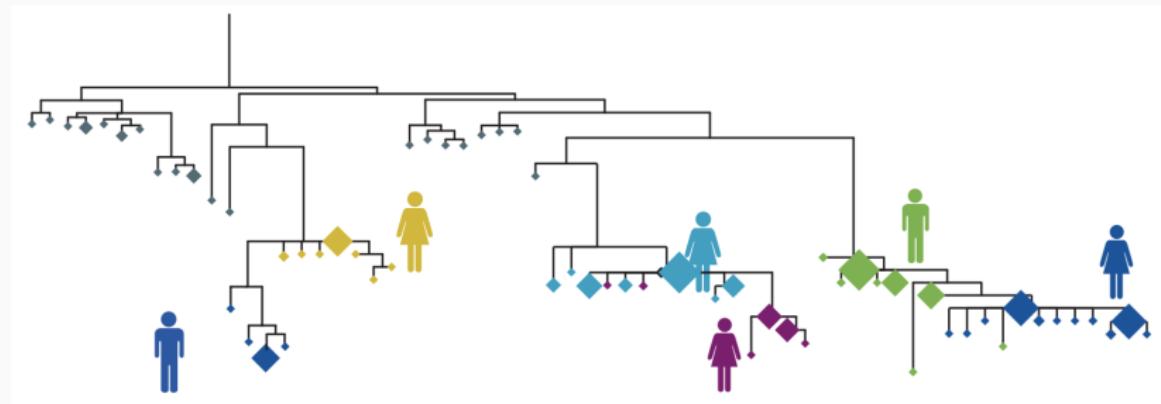


Figure from Cluster Picker (<https://hiv.bio.ed.ac.uk/software.html>)

- Cluster analysis → how many, how big, importations
- **Ancestral state reconstruction**
 - Characterising epidemiological transmission dynamics → source attribution
 - Phylogeography → spatial transmission dynamics
- Phylodynamics → estimating population-level parameters which shape phylogenies

Ancestral state reconstruction

Identifying likely transmission pairs for source attribution



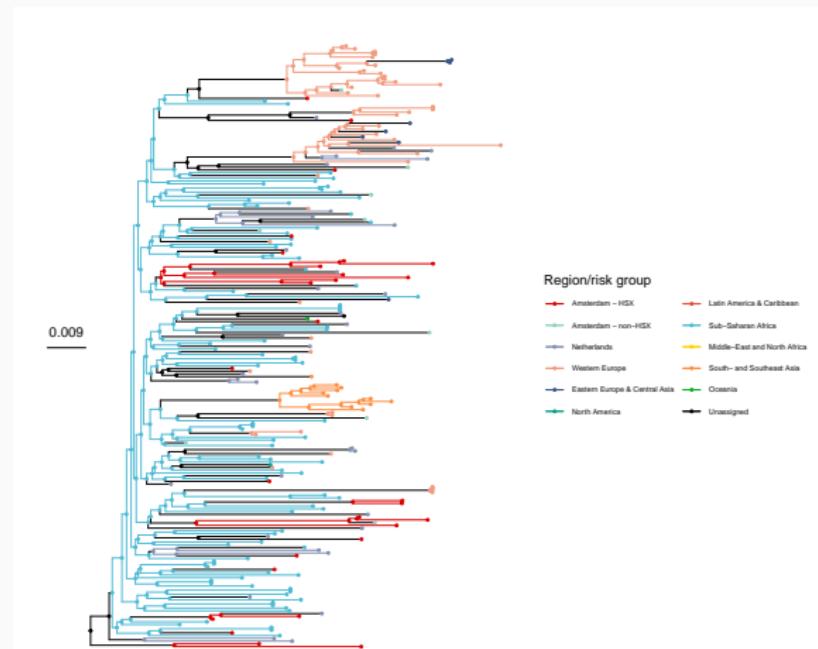
Ratmann et al (2020), *Lancet HIV*

Ancestral state reconstruction

Separating phylogenetic transmission chains in study population

Software:

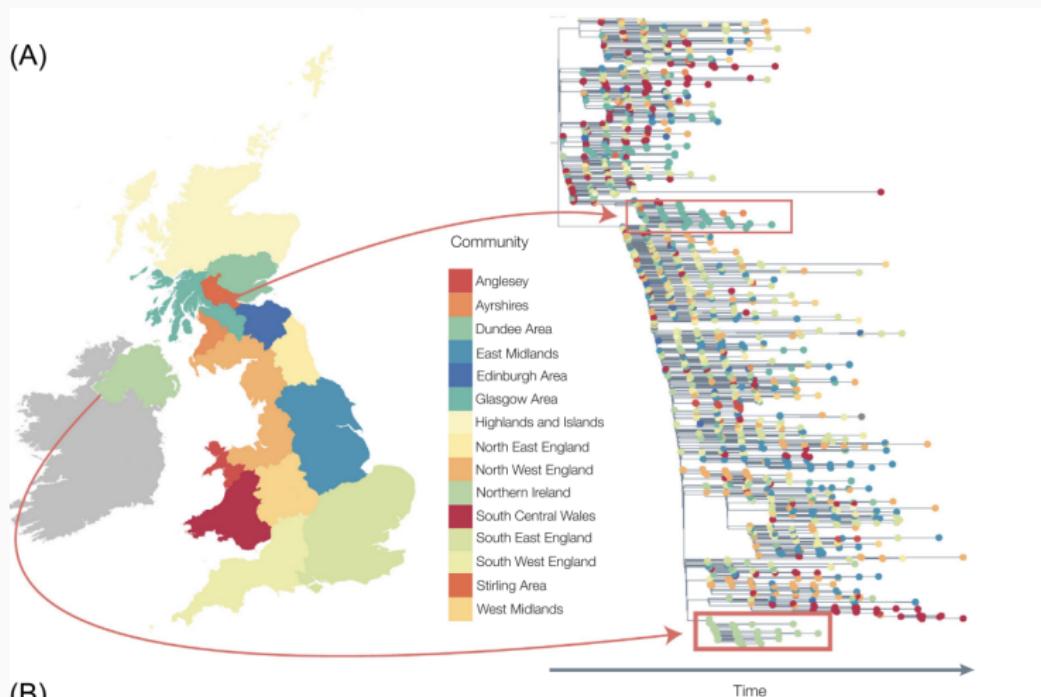
- *Phyloscanner*
- *BEAST*



Blenkinsop et al (2022), *eLife*

Ancestral state reconstruction

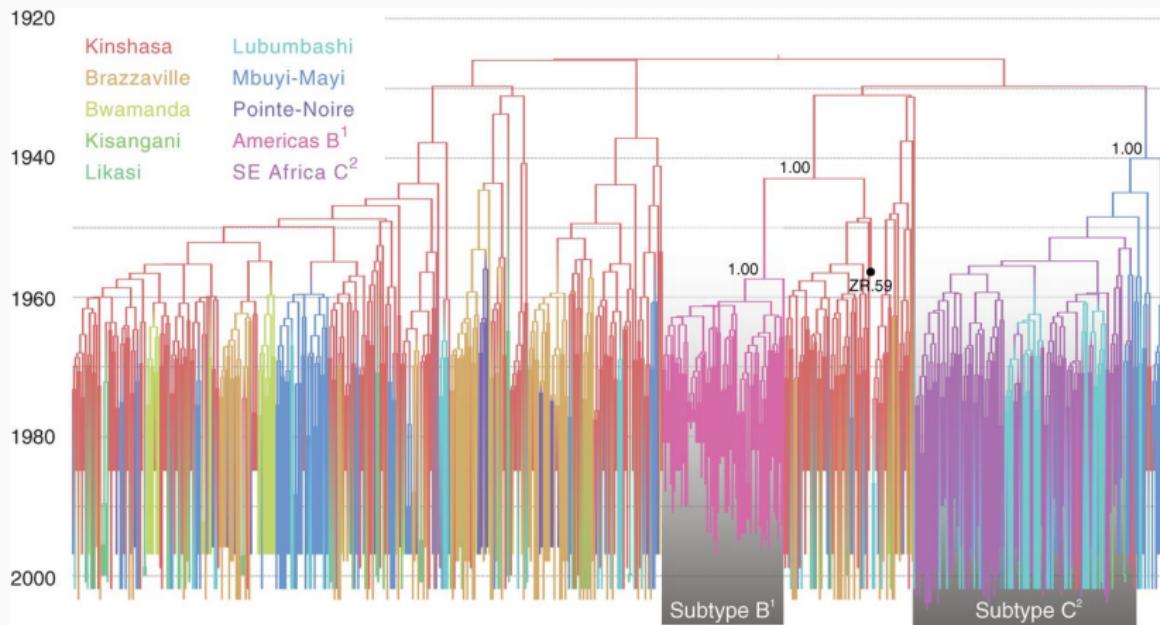
Tracking outbreaks across spatial domains



Hill et al (2021), *Trends in Parasitology*

Ancestral state reconstruction

Inferring the spatial epidemiological history of a virus

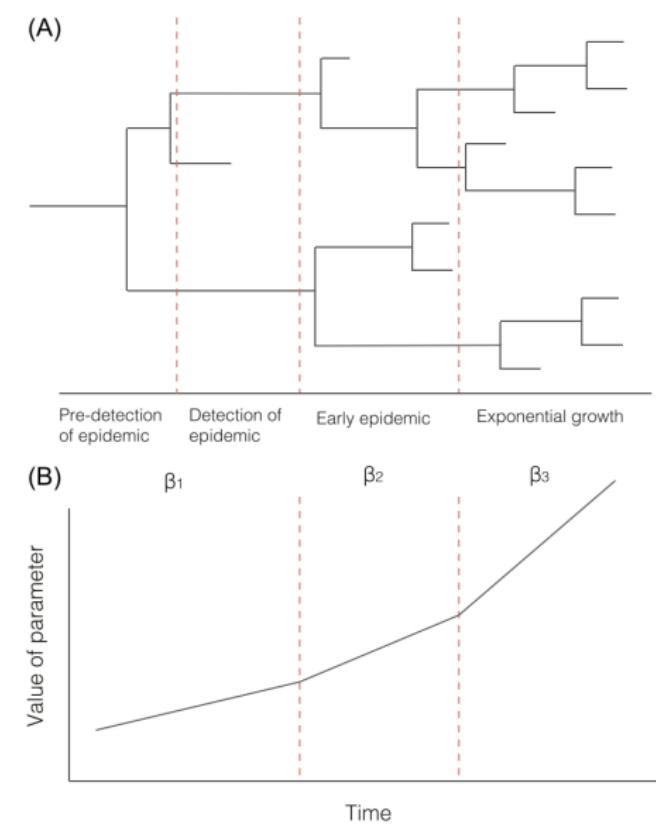


Faria et al. (2014), *Science*

- Cluster analysis → how many, how big, importations
- Ancestral state reconstruction
 - Characterising epidemiological transmission dynamics → source attribution
 - Phylogeography → spatial transmission dynamics
- **Phylodynamics → estimating population-level parameters which shape phylogenies**

Viral phydynamics

Characterising underlying processes which shape viral phylogenies
- linking epidemiology with evolutionary dynamics



- Genomic data can provide insights into dynamics of virus spread at different scales
- Phylogenetic trees are a natural way to describe ancestry
- Reconstructing ancestral states helps to understand epidemiological patterns
- Utilising phylogenetic data in statistical epidemiological models enable us to make population-level inferences

Acknowledgements

Thanks to Chris Wymant, Thibaut Jombart and PANGEA, whose slides I adapted.

