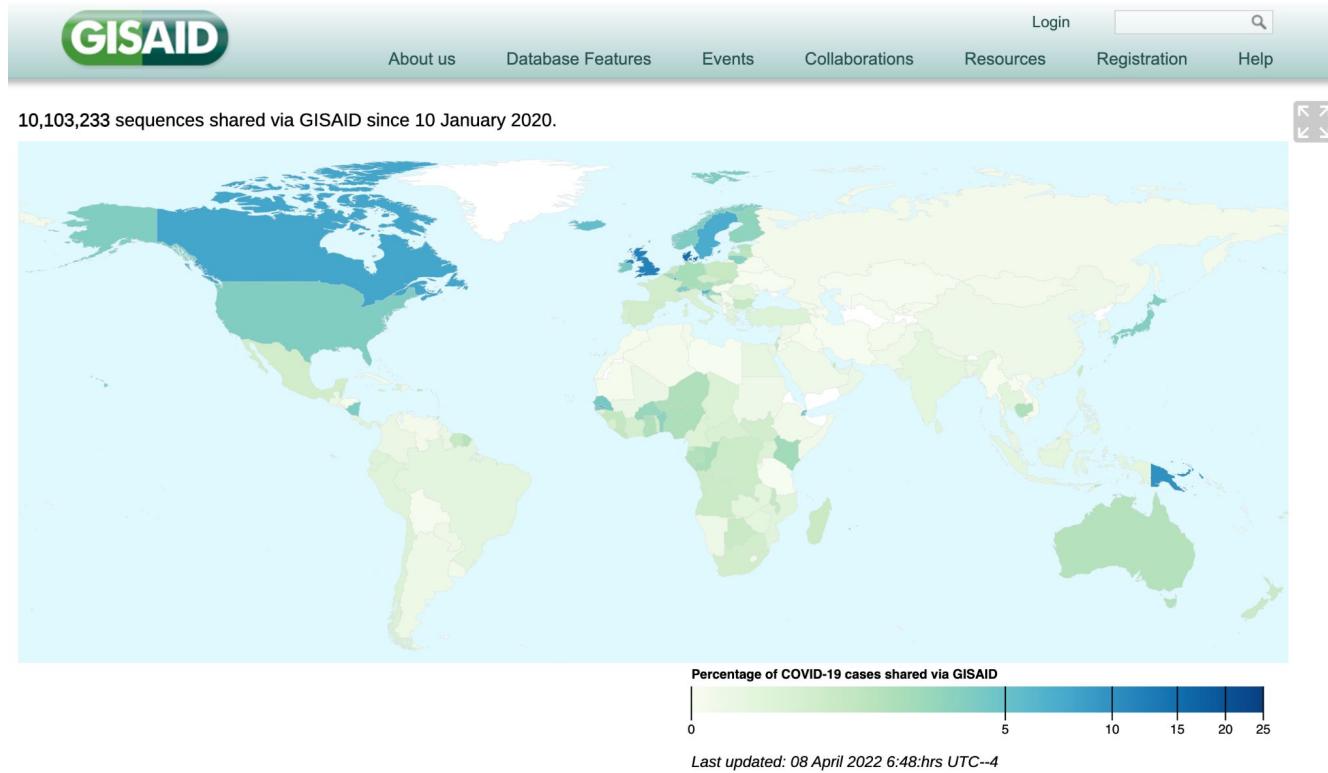


After the data deluge: scaling strategies for massive genomic datasets

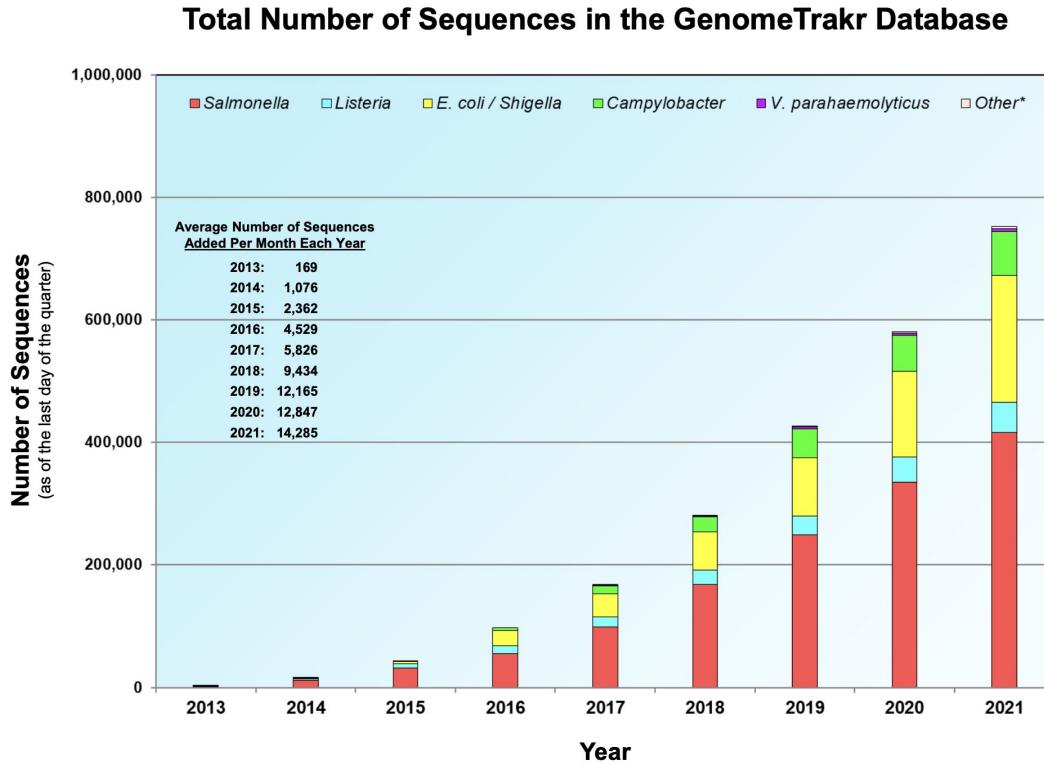
Molecular Epidemiology of Infectious Diseases
Lecture 11

April 11th, 2022

The data deluge



The data deluge



**Massive genomic
datasets have pushed
the limits of existing
methods**

The benefits of Bayesian inference have a cost

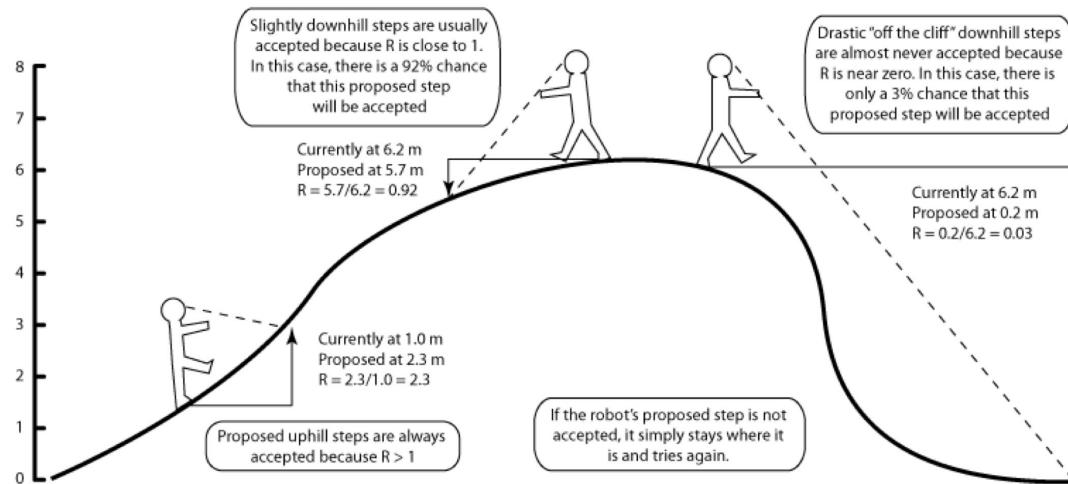
We have primarily focused on Bayesian phylogenetic methods because they allow for flexible, probabilistic modeling of pathogen evolution

While quantifying our uncertainty in epi/evo parameters as well as the pathogen phylogeny.

BUT these Bayesian methods are highly reliant on MCMC sampling.

The costs of Bayesian MCMC

Like a blind robot on a random walk, MCMC is inherently inefficient



Problems with MCMC efficiency

Sampling via a random walk is inherently inefficient as it may take a long time for the chain to converge on the posterior distribution.

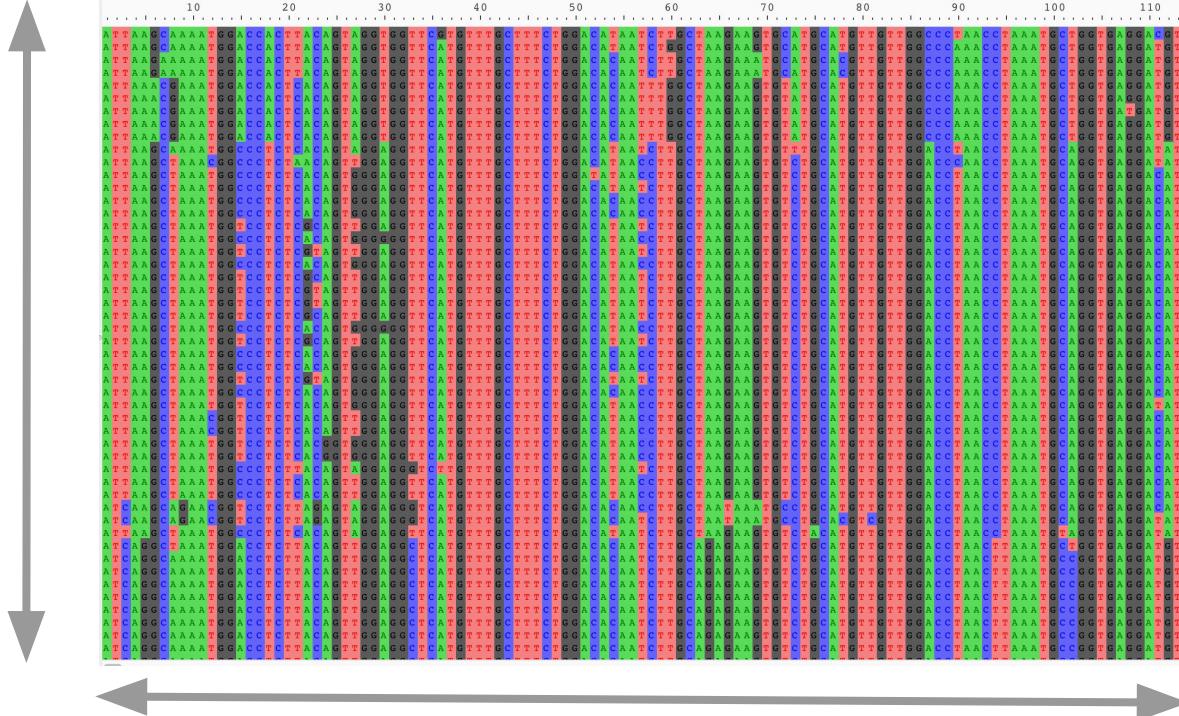
Using a Metropolis-Hastings step means many proposals are rejected.

Samples are generally highly autocorrelated, requiring a lot of wasted computation (thinning) to get pseudo-independent samples from the posterior.

Faster gradient-based MCMC methods (e.g. Hamiltonian Monte Carlo, NUTS) cannot be easily adapted to sampling from tree space.

What do we even mean by big?

Long (many samples)



Wide (many sites/characters)

**What is the bigger
problem: the length
or width of a genomic
dataset?**

Scaling with the number of sites

The likelihood of the sequence data is generally computed independently at each site.

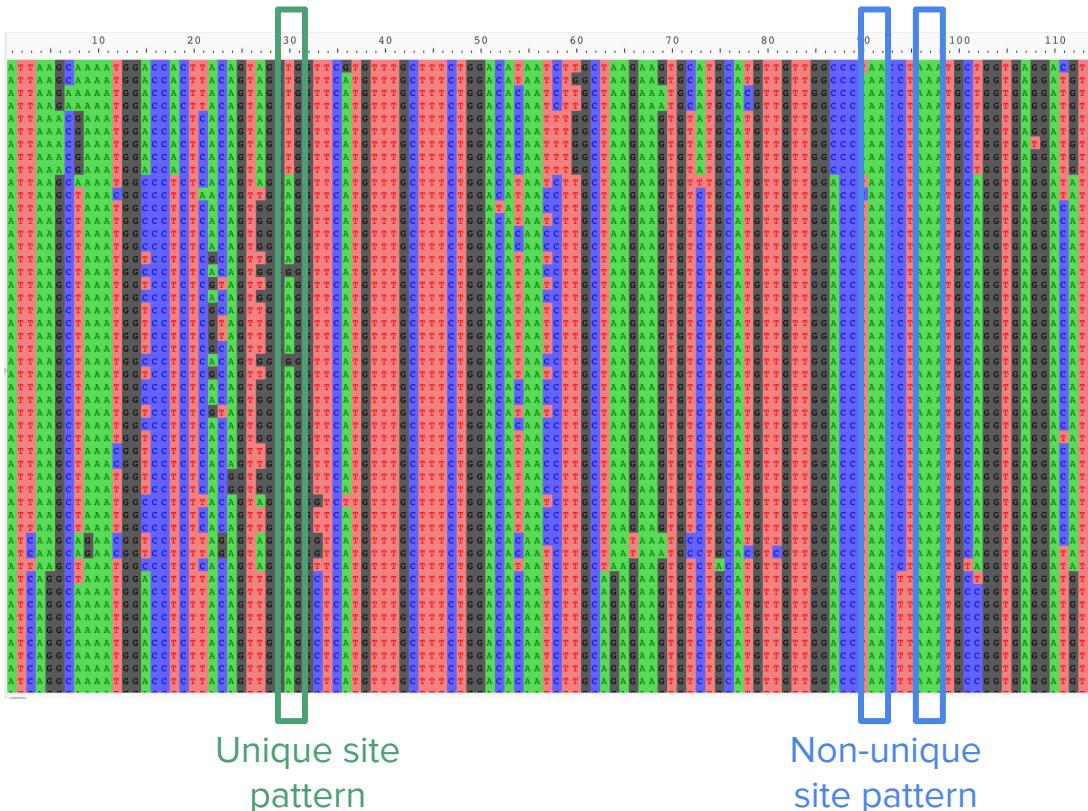
$$L(Seq|Tree) = \prod_{i=1}^{i=N} L(Seq_i|Tree)$$

Computation time therefore generally increases linearly with the number of sites (width) of an alignment.

Or even less than linearly because we actually only need to compute the likelihood once for each **unique site pattern** in an alignment.

Site patterns

Many site patterns will not be unique such that the cost of computing the likelihood will generally increase slightly less than linearly with the number of sites in an alignment.

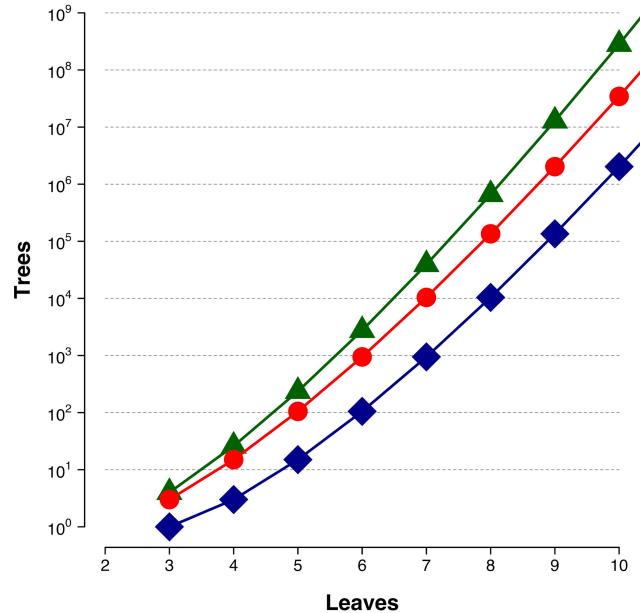


Scaling with the number of samples

On the other hand, tree space grows exponentially with the number of tips or samples

We also need to estimate an additional branch length parameter for each sample we add.

MCMC methods thus scale very poorly with the number of samples

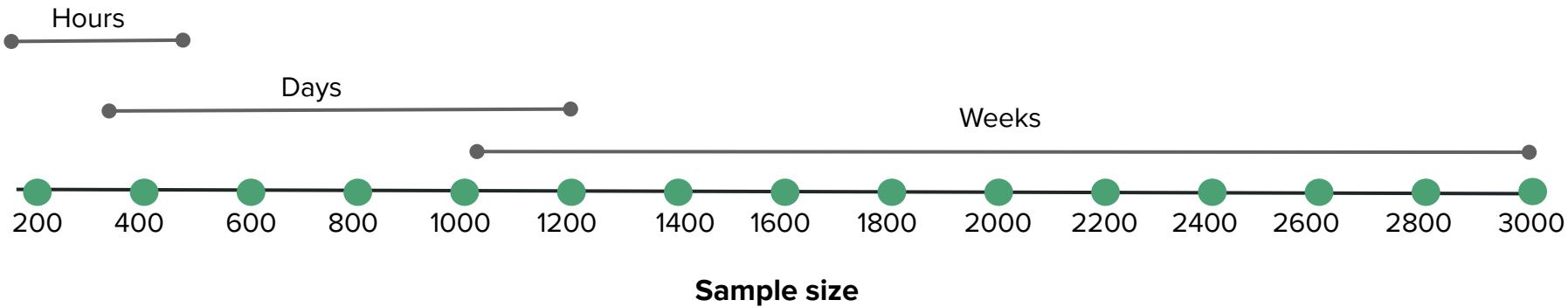


Red shows rooted binary trees.

Practical limits of Bayesian phylogenetics

Generally Bayesian phylogenetic inference using MCMC is limited to at most a few thousand samples.

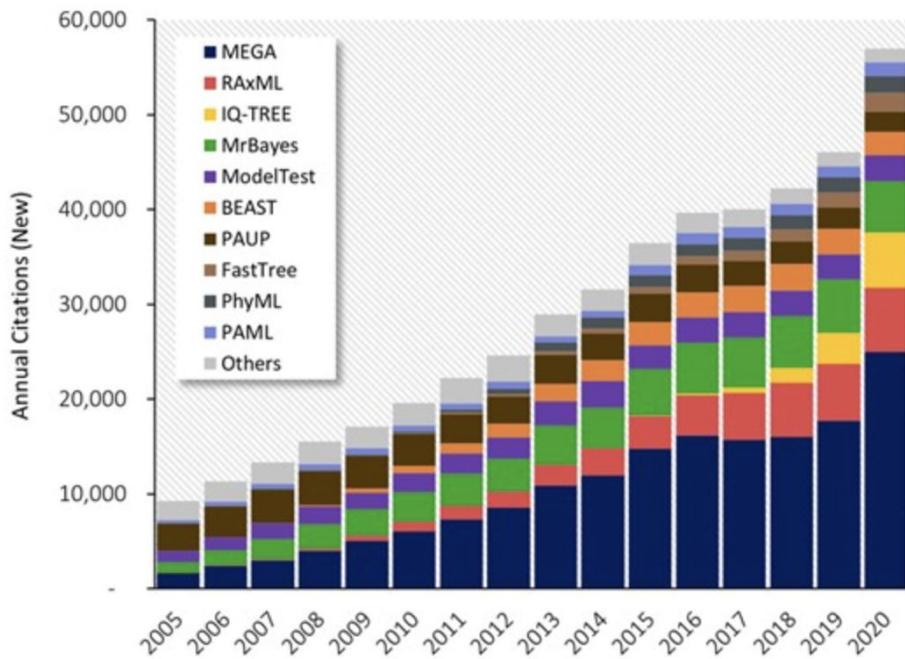
Assuming a typical 10kb alignment, the figure below gives a rough approximation of how run times scale with sample size.



Faster tree reconstruction methods

Nothing wrong with ML

A large number of biologists still use ML or other fast tree reconstruction methods.



Kumar (MBE, 2022)

Fast ML phylogenetic methods

Many popular ML-based methods use a similar reconstruction strategy:

1. A starting tree is obtained using a very fast approximate method like Neighbor Joining or Maximum Parsimony.
2. A hill climbing algorithm is used to search for trees with higher likelihood:
 - a. A new tree is proposed by a tree-rearrangement move like NNI or SPR.
 - b. Tree is accepted only if it increases the likelihood of the sequence data.
3. The search stops once no further improvements to the likelihood can be found
4. Optional: The entire process starting from step 1 is repeated with a different starting tree.

Fast ML phylogenetic methods

Zhou *et. al.* compared the performance of four of the most popular ML methods.

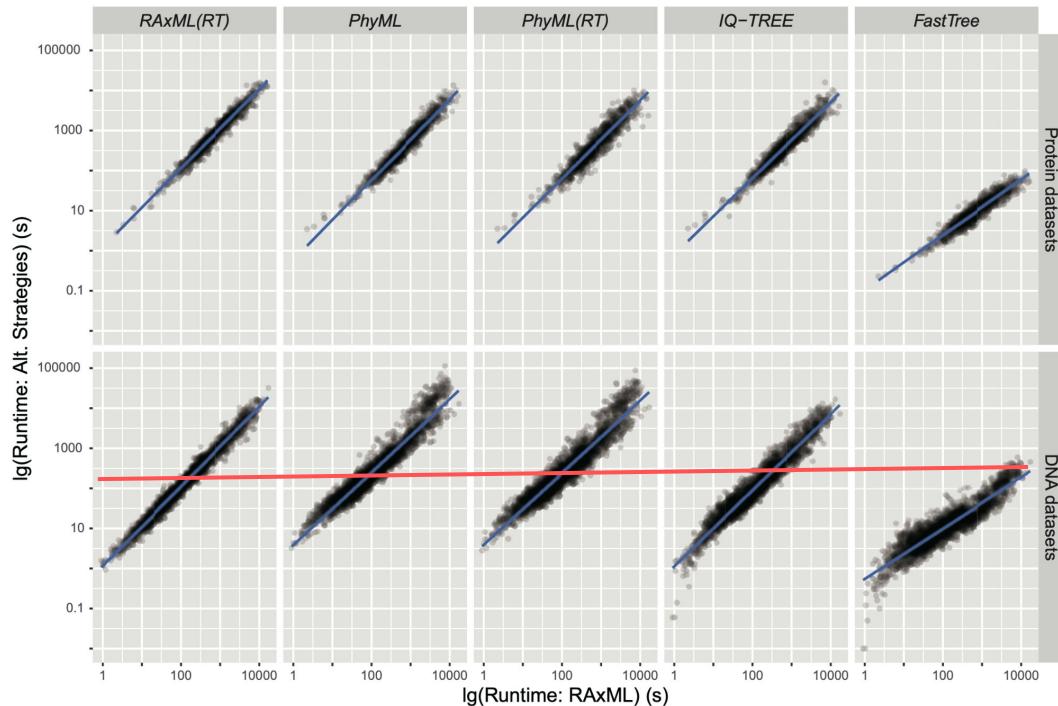
Table 1. Overview of the Four Fast ML-Based Phylogenetic Programs Evaluated in This Study.

Programs	Optimality Criterion	Starting Tree	Topological Moves	Supported Models		Partitioned Analysis
				AA	DNA	
RAxML v8.2.0 (ExaML v3.0.17)	ML	Parsimony/random/custom	SPR	Common and custom models	JC69, K80, HKY85, GTR	Y
PhyML v20160530	ML	Parsimony/random/custom	Interleaved NNI and SPR	Common and custom models	Common and custom models	Y
IQ-TREE v1.4.2	ML	BIONJ and multiple parsimony/random/custom	NNI and stochastic perturbation	Common and custom models	Common and custom models	Y
FastTree v2.1.9	ML	Heuristic NJ	NNI and SPR (ME) followed by NNI (ML)	JTT, WAG, LG	JC69, GTR	N

NOTE.—ML, maximum likelihood; ME, minimum evolution; NJ, neighbor joining; NNI, nearest neighbor interchange; SPR, subtree pruning and re-grafting.

Fast ML phylogenetic methods

FastTree can be 10-100X times faster than similar ML tree approaches.

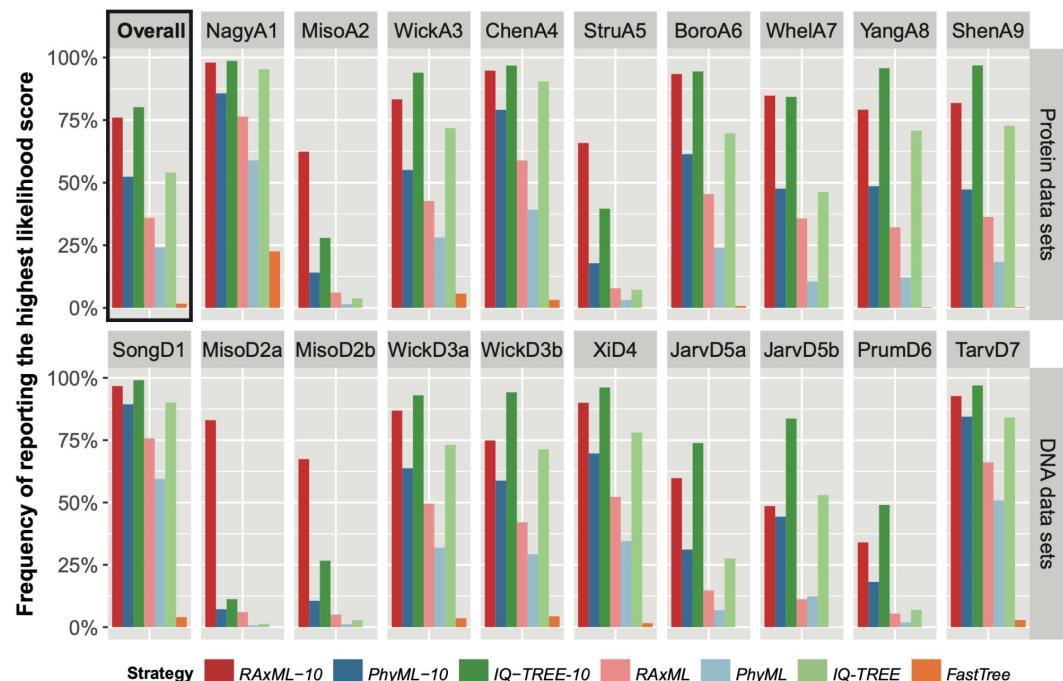


Fast ML phylogenetic methods

RAxML, PhyML and IQ-TREE perform about equally well.

Using multiple starting trees (x-10 runs) considerably helps in finding the best tree.

FastTree is considerably less accurate and rarely recovers the highest likelihood tree.



Fast ML phylogenetic methods

There is an essential tradeoff between speed and accuracy determined by how exhaustively different ML methods search tree space.

My vote: These days I generally use RAxML or IQ-TREE. IQ-TREE might be gaining an edge since it has some nice features like built-in dating and ultra-fast bootstrap searches.

But I still use FastTree for explorative analyses with very large alignments (>10,000 sequences) when fine-scale accuracy is less important than speed.

Don't “reinvent the tree”

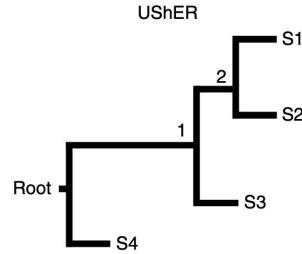
During epidemics, the number of sequences will generally grow rapidly over time, requiring *online* or *real-time* methods that can quickly incorporate new data.

Many fast ML methods like IQ-TREE and RAxML (via EPA-ng) allow for new samples to be placed on existing trees. This is known as ***sample placement***.

This can however still be time consuming. Turakhia *et al.* (2021) found that it takes about 28 minutes to place one SARS-CoV-2 sample on an existing reference tree with 38,342 tips.

Ultra-fast placement algorithms

UShER uses a maximum parsimony approach to search a *mutation-annotated tree* for a placement that requires the fewest additional mutations.

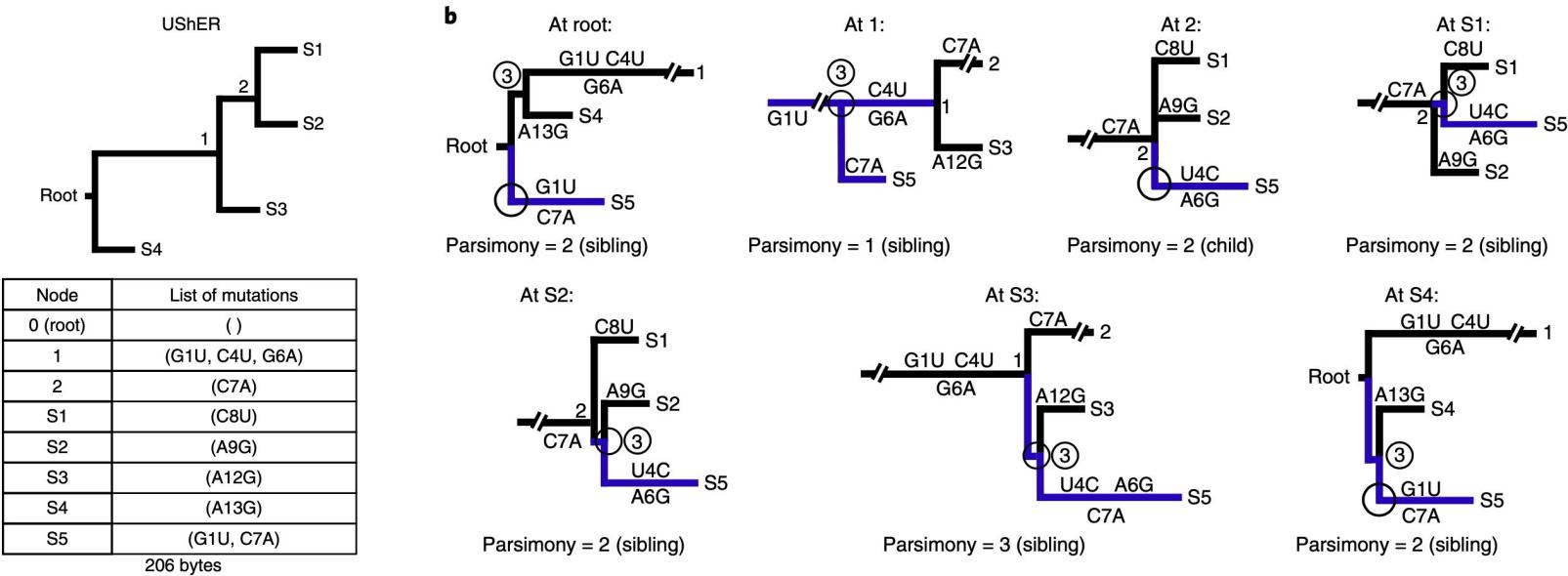


Node	List of mutations
0 (root)	()
1	(G1U, C4U, G6A)
2	(C7A)
S1	(C8U)
S2	(A9G)
S3	(A12G)
S4	(A13G)
S5	(G1U, C7A)

206 bytes

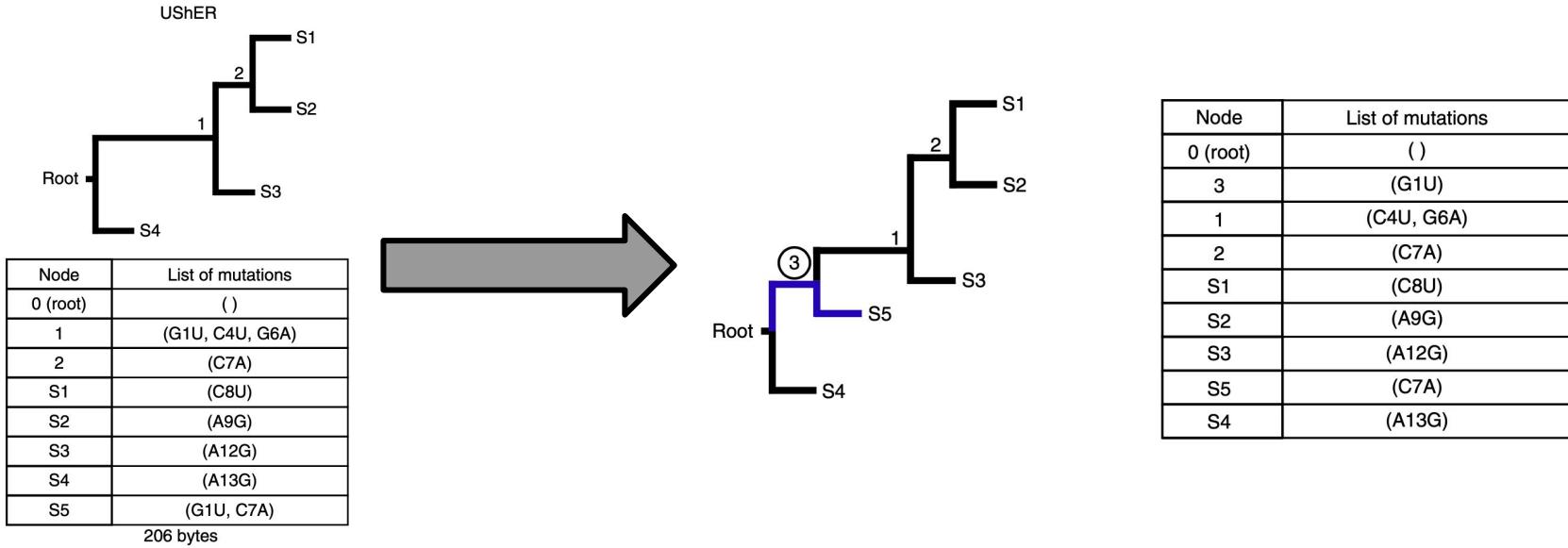
Ultra-fast placement algorithms

UShER uses a maximum parsimony approach to search a *mutation-annotated tree* for a placement that requires the fewest additional mutations.



Ultra-fast placement algorithms

UShER uses a maximum parsimony approach to search a *mutation-annotated tree* for a placement that requires the fewest additional mutations.



Ultra-fast placement algorithms

Fast placement algorithms allow for ‘real-time’ phylogenetic inference during large epidemics by adding new samples and only periodically rebuilding the reference tree.

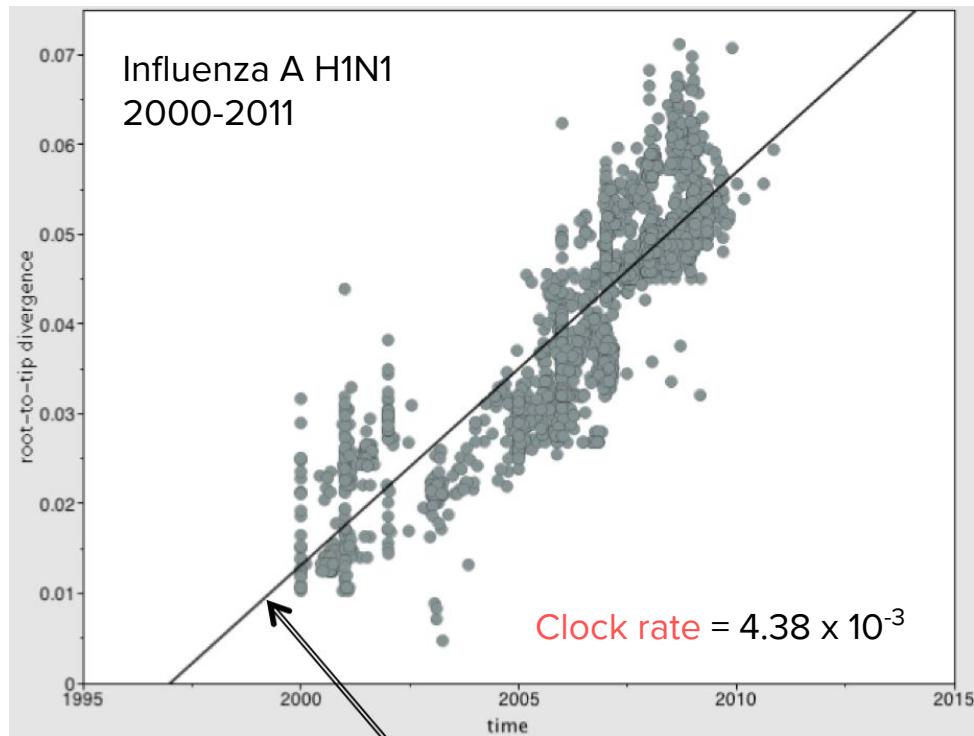
Works well when genetic diversity is limited and sequences are only separated by a few mutations.

Danger zone: Parsimony is known to be less accurate than ML when there’s a reasonably high probability that multiple mutations (including reversions) have occurred along a branch.

Faster dating: methods for time-calibrated phylogenies

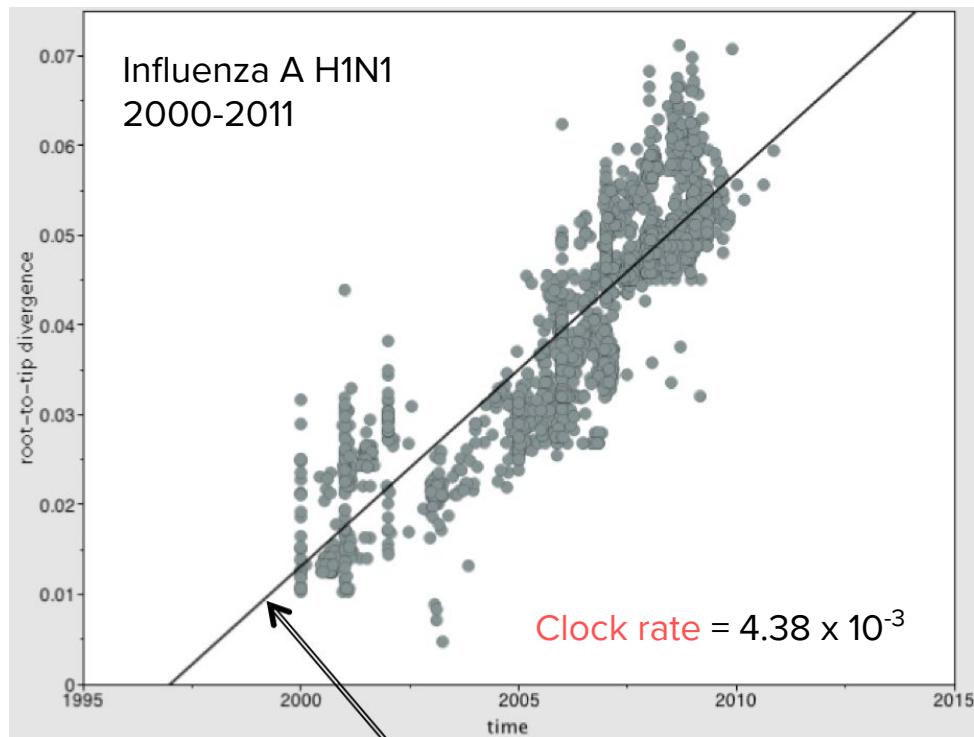
Root-to-tip regression

Recall that under a strict molecular clock model lineages will accumulate mutations at a constant rate such that the number of mutations between a tip and the root should increase linearly with time.



Root-to-tip regression

Thus if we regress root-to-tip divergence (in number of mutations) against sampling times, the slope of the resulting regression line will give us a rough estimate of the molecular clock rate.



Least-squares dating

Least-squares dating methods like LSD (To *et al.*, 2016) use root-to-tip regression to estimate clock rates and node heights (i.e. divergence times).

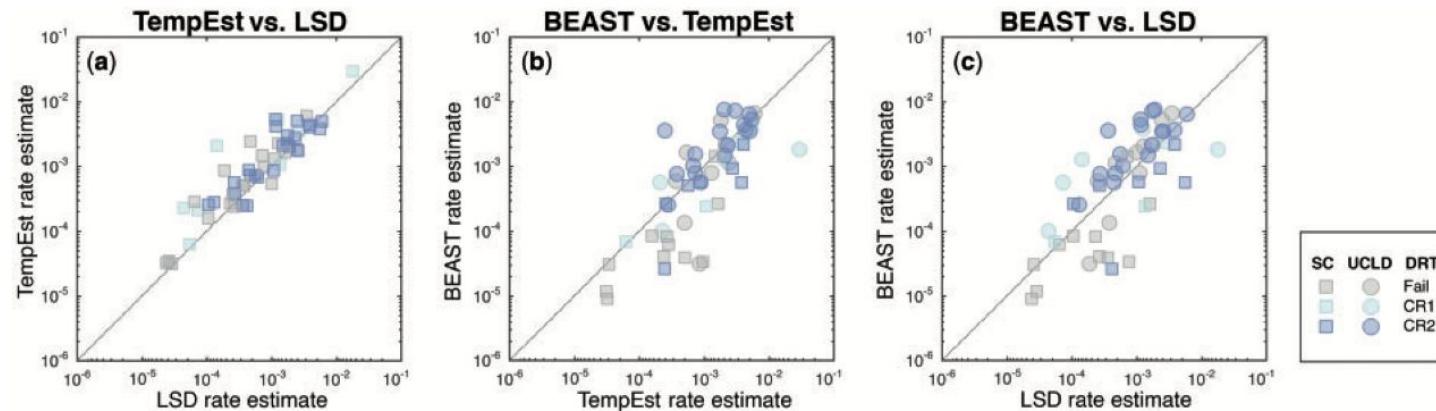
The node heights t_i and clock rate ω can be estimated together using least-squares optimization:

$$\phi(\omega, t_1, \dots, t_{n-1}) = \sum_{i=2}^{2n-1} \frac{1}{\sigma_i^2} (b_i - \omega(t_i - t_{a(i)}))^2$$

This is very similar to standard least-squares regression: we try to minimize the deviation between the observed branch lengths b_i and the expected number of mutations $\omega(t_i - t_{a(i)})$ along that branch under the molecular clock assumption.

How do dating methods compare?

Regression-based methods like LSD generally return molecular clock estimates that are very similar to Bayesian dating methods.



Other fast dating methods

TempEst: Simple GUI app for root-to-tip regression from the creators of BEAST (Rambaut *et al.*, 2016). See tutorial from week one for an example.

TimeTree: Iteratively optimizes branch lengths and ancestral sequences in a maximum likelihood framework. Available in the *TimeTree* Python package (Sagulenko *et al.*, 2018).

BactDating: Allows for Bayesian inference of divergence times from a fixed tree. Quantifies uncertainty in divergence times and partially accounts for recombination (Didelot *et al.*, 2018).

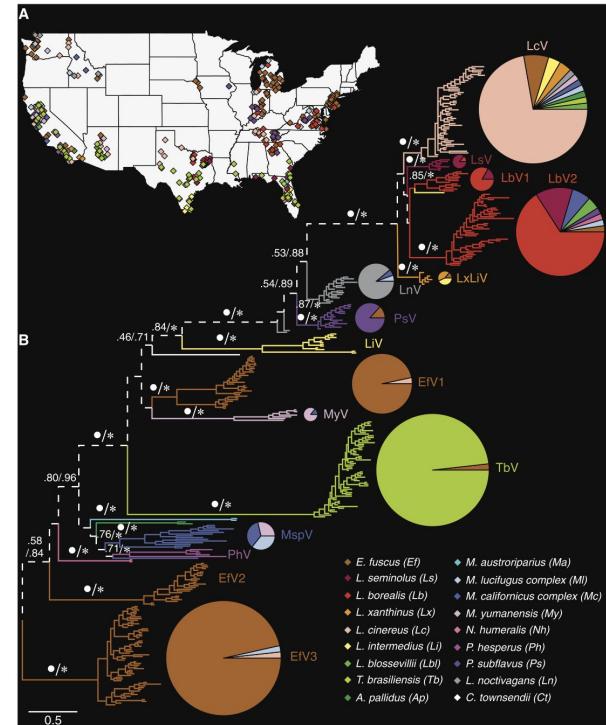
Faster ancestral state reconstruction

Ancestral state reconstruction

Ancestral state reconstruction is widely used to reconstruct changes in the state of lineages along a tree.

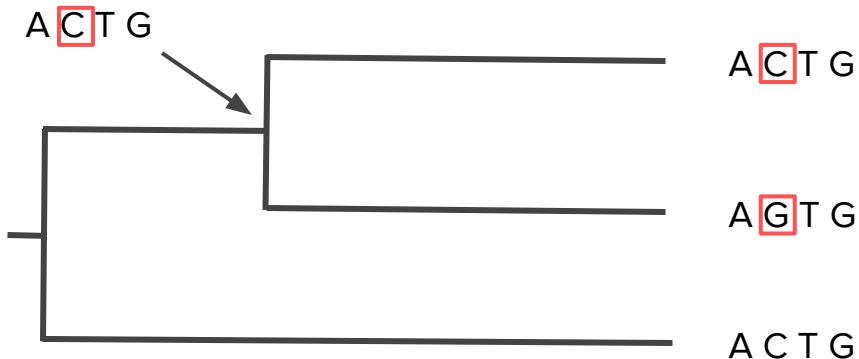
In phylogeography, ancestral states are often locations but a lineage's state could refer to its host population, genotype or even its entire sequence.

Given a fixed tree, maximum parsimony (MP) and maximum likelihood (ML) methods allow for very efficient reconstructions.



Streicker et al. (Science, 2010)

Likelihood of sequence data on trees



Maximum likelihood methods work very similar to how we computed the likelihood of sequence data given a tree using Markovian evolutionary models.

Modeling molecular evolution

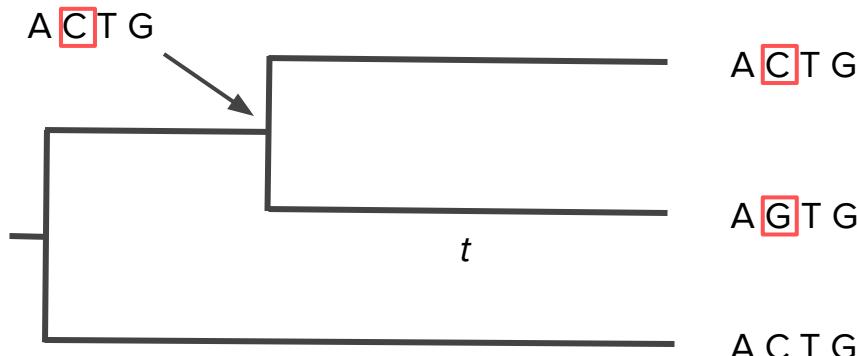
We can compute transition probabilities under a continuous-time Markov model given our substitution matrix \mathbf{Q} and the time elapsed along a branch t .

$$P(t) = e^{Qt}$$

The elements of $P(t)$ give us the probability of every possible transition. Importantly, these **transition probabilities take into account every possible substitution path**.

$$P(t) = \begin{bmatrix} P_{T,T} & P_{T,C} & P_{T,A} & P_{T,G} \\ P_{C,T} & P_{C,C} & P_{C,A} & P_{C,G} \\ P_{A,T} & P_{A,C} & P_{A,A} & P_{A,G} \\ P_{G,T} & P_{G,C} & P_{G,A} & P_{G,G} \end{bmatrix}$$

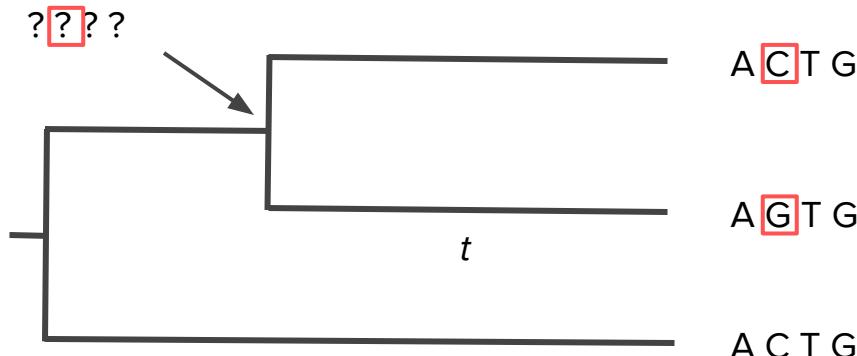
Computing likelihoods at one site



We can compute the likelihood of the sequence data given each possible ancestral state:

$$L(Seq|Tree) = P_{C,C}(t) * P_{C,G}(t)$$

Computing likelihoods at one site



Then find the ancestral state that maximizes the likelihood of the sequence data:

$$\text{ML state} = \underset{X \in \{A,C,T,G\}}{\operatorname{argmax}} (P_{X,C}(t) * P_{X,G}(t))$$

ML ancestral state reconstruction

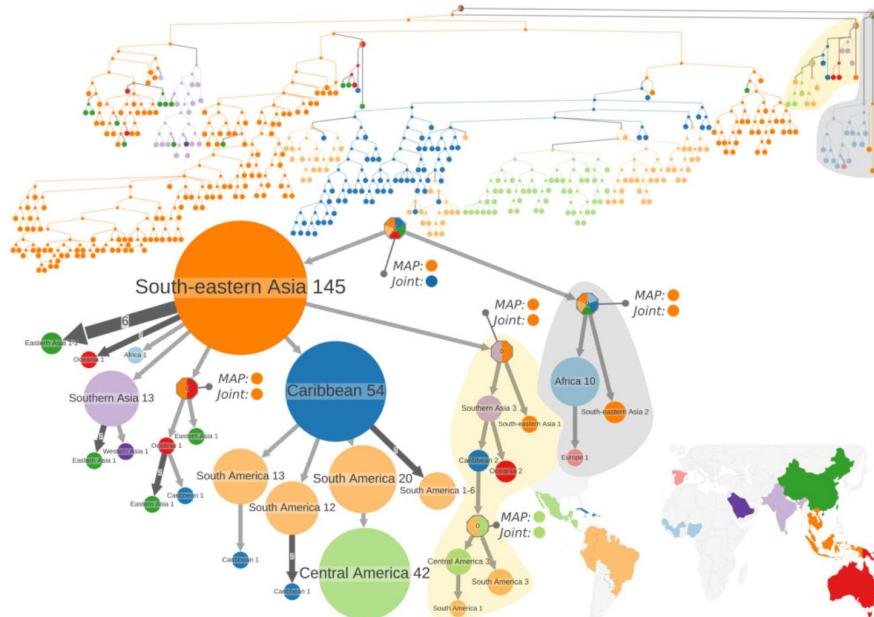
ML inference allows for state reconstruction under standard Markovian models of sequence or character evolution like the General-Time-Reversible (GTR) model that allow for transition rates to differ between states.

While ML inference does not allow for full quantification of uncertainty about all parameters, it allows for ancestral state probabilities to be computed using the relative likelihood of each state.

Methods like PastML (Ishikawa *et al.*, 2019) allow for efficient inference of both transition rates and ancestral states.

Ancestral reconstructions with PastML

PastML allows for ancestral histories to be compressed into simpler graphs like in this example showing the global dissemination of dengue virus serotype 2.



MP ancestral state reconstruction

Maximum parsimony reconstructions assume a “minimal evolution” model that minimizes the number of state changes required to explain the observed states at the tips of the tree.

MP works well when transition rates are low such that the probability of multiple state changes along a lineage is low. In this case, the most likely ancestral state mapping will also be the most parsimonious.

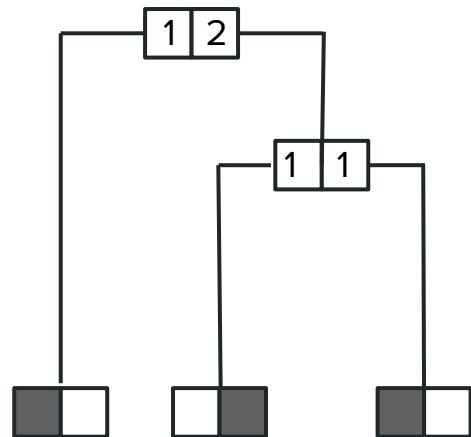
But with more rapid character evolution, the most parsimonious history may not be the most likely as it may be unlikely that only a single state change has occurred.

MP algorithms

MP allows for ultra-fast reconstructions using clever dynamic programming methods.

Generally these involve just two steps:

A **post-order** (tip-to-root) traversal: The number of changes required to explain the tip data subtending each internal node are computed



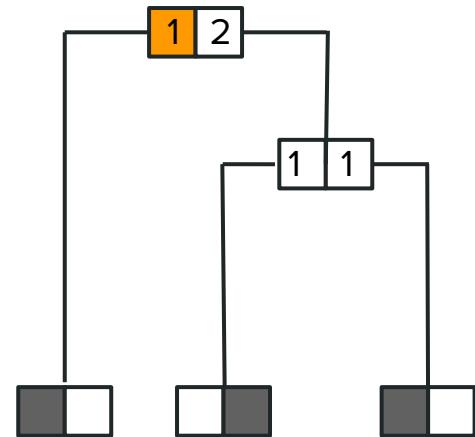
MP algorithms

MP allows for ultra-fast reconstructions using clever dynamic programming methods.

Generally these involve just two steps:

A ***post-order*** (tip-to-root) traversal: The number of changes required to explain the tip data subtending each internal node are computed

A ***pre-order*** (root-to-tip) traversal: An ancestral node is chosen at each internal node that minimizes the number of changes required above and below the node.



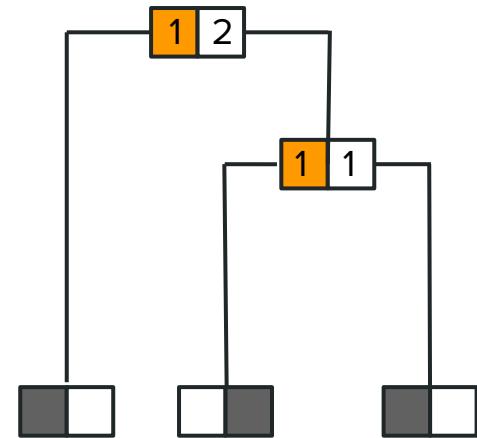
MP algorithms

MP allows for ultra-fast reconstructions using clever dynamic programming methods.

Generally these involve just two steps:

A ***post-order*** (tip-to-root) traversal: The number of changes required to explain the tip data subtending each internal node are computed

A ***pre-order*** (root-to-tip) traversal: An ancestral node is chosen at each internal node that minimizes the number of changes required above and below the node.



MP ancestral state reconstruction

The main advantage of MP reconstructions is speed. There are no parameters to estimate and reconstructions generally take less than a second even for trees with thousands of tips.

Thus ultra-fast placement methods like UShER also use MP to place additional samples on existing trees.

Popular phylogenetic software packages like MEGA and Mesquite implement MP reconstructions.

Learning to let go: sub- and re-sampling strategies

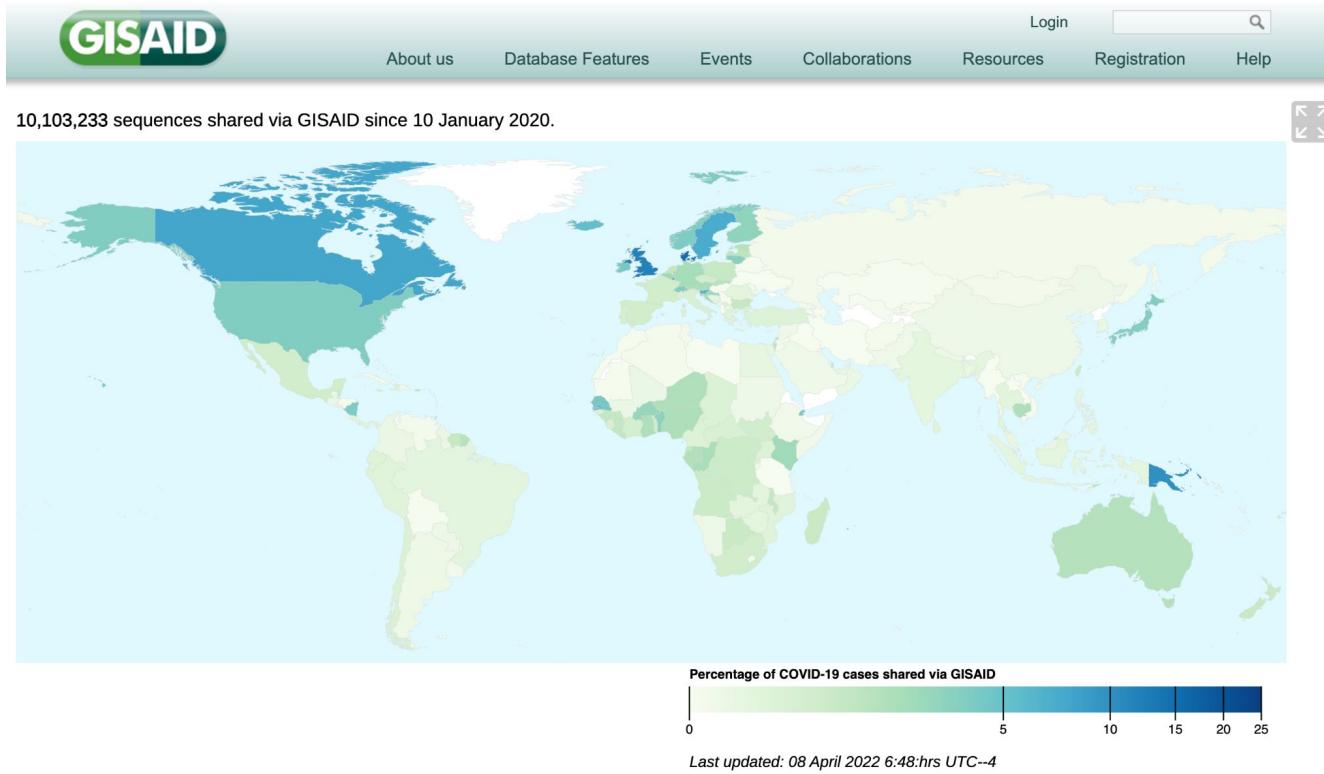
Subsampling

It's important to remember that it is often possible to answer the questions of interest to us with much smaller datasets than the total available data.

Subsampling strategies allow us to down-sample data sets to more manageable sizes while correcting for biased sampling in the original dataset.

Resampling strategies: allow us to explore how sampling choices may be influencing our inferences.

Global bias in SAR-CoV-2 sampling



Reconstructing CoV imports and exports

RESEARCH

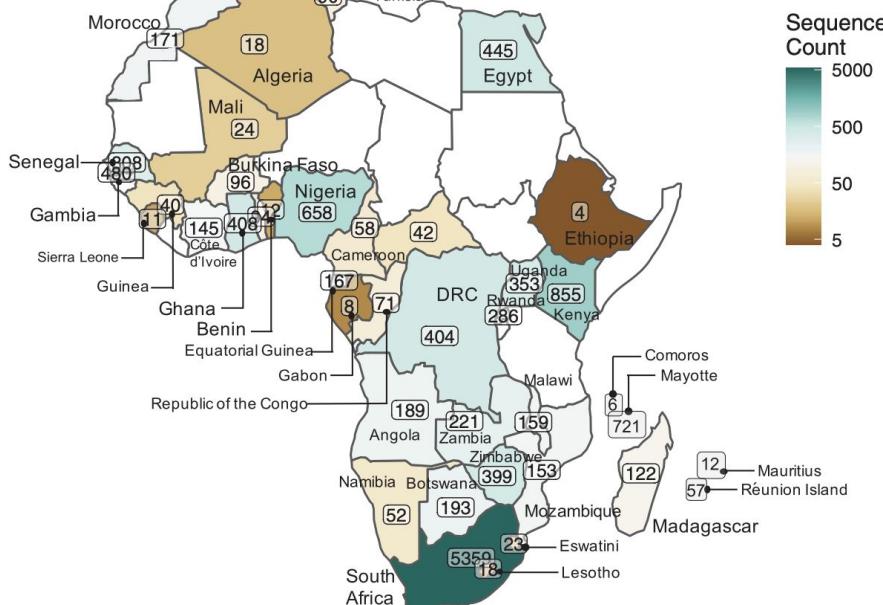
RESEARCH ARTICLE

CORONAVIRUS

A year of genomic surveillance reveals how the SARS-CoV-2 pandemic unfolded in Africa

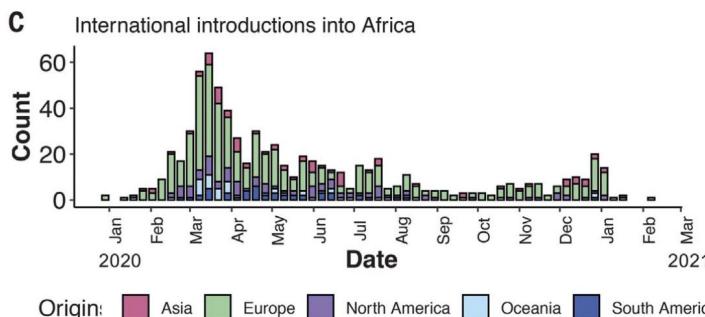
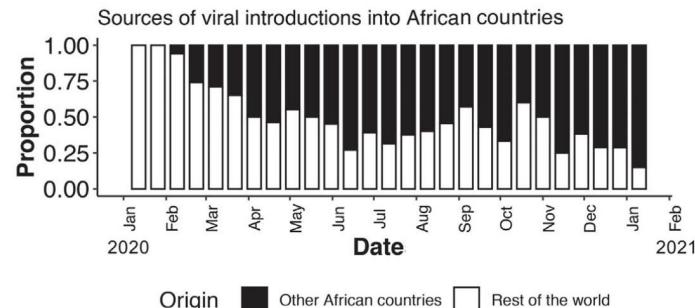
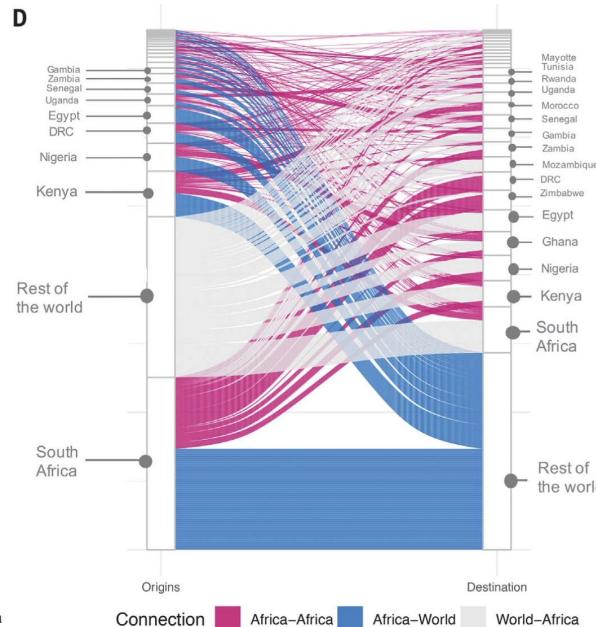
The progression of the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) pandemic in Africa has so far been heterogeneous, and the full impact is not yet well understood. In this study, we describe the genomic epidemiology using a dataset of 8746 genome from 33 African countries and two overseas territories. We show that the epidemic in most countries was initiated by importations predominantly from Europe, which diminished after the early introduction of international travel restrictions. As the pandemic progressed, ongoing transmission in many countries and increasing mobility led to the emergence and spread within the continent of many variants of concern and interest, such as B.1.351, B.1.925, A.231, and C.1. Although controlled by low sampling numbers and blind spots, the findings highlight that Africa must not be left behind in the global pandemic response, otherwise it could become a source for new variants.

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) emerged in late 2019 in Wuhan, China (*1,2*). Since then, the virus has spread to all corners of the world, causing almost 150 million cases of COVID-19 and more than 3 million deaths by the end of April 2021. Throughout the pandemic, it has been noted that Africa accounts for a relatively low proportion of reported cases and deaths—by the end of April 2021, there had been ~4.5 million cases and ~120,000 deaths on the continent, corresponding to less than 4% of the global burden. However, emerging data from seroprevalence surveys and autopsy studies in some African countries suggest that the true number of infections and deaths may be several-fold higher than reported (*3,4*). In addition, a recent analysis showed that



Reconstructing CoV imports and exports

A large number of imports were identified early in the epidemic, mostly from Europe.



Reconstructing CoV imports and exports

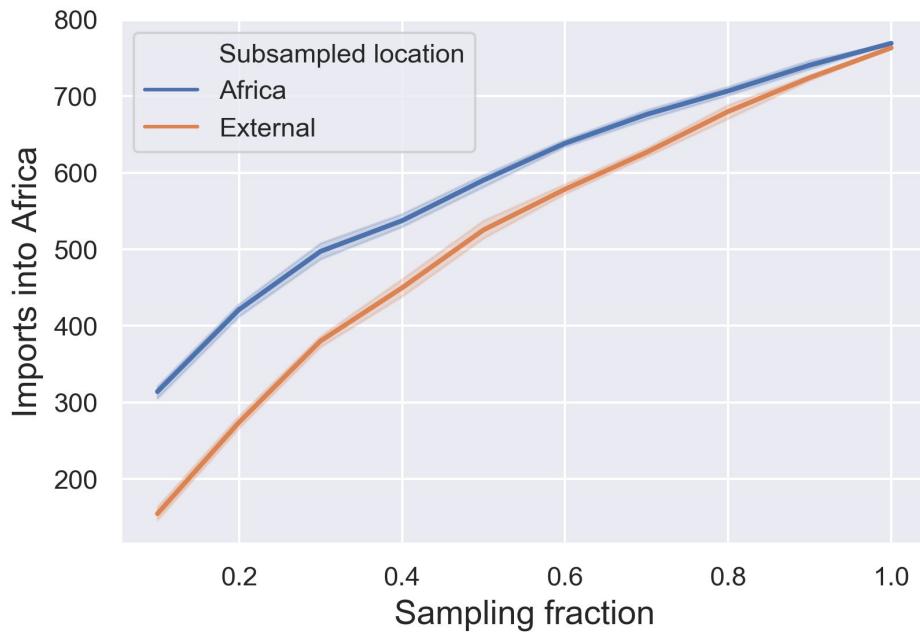
However, the number of imports/exports identified will almost certainly depend on how pathogen genomes were sampled in and outside of Africa.

We therefore performed a rarefaction analysis where we systematically varied the fraction of sampled genomes from either Africa or the rest of the world.

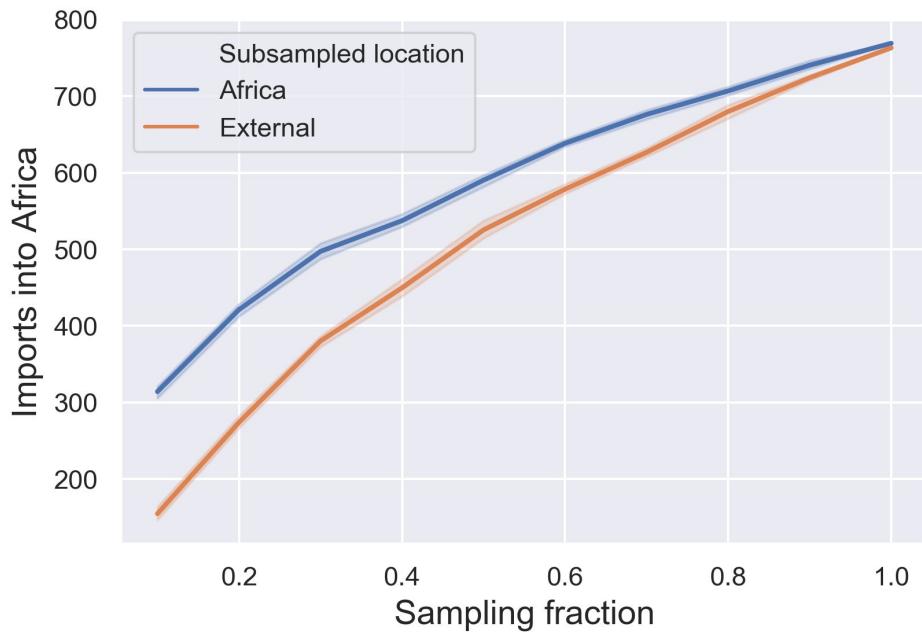
Basic strategy:

1. Build ML tree for full data set
2. Subsample a fraction x of genomes for a given location
3. Reconstruct ancestral locations using MP based on new subsampled tree
4. Count number of imports/exports identified at sampling fraction x

Reconstructing CoV imports and exports



Reconstructing CoV imports and exports



Vastly more introductions would have been identified with increased sampling in Africa or globally, suggesting that the intros identified are really just the “ears of the hippo”.



Final thoughts

After the data deluge



Soldiers disinfecting parts of Brasilia's underground rail network as the coronavirus spread throughout Brazil in late March 2020.

<https://www.nature.com/articles/d41586-021-00525-x>

**Want to track pandemic variants faster?
Fix the bioinformatics bottleneck**

Emma B. Hodcroft, Nicola De Maio, Rob Lanfear, Duncan R. MacCannell, Bui Quang Minh,
Heiko A. Schmidt, Alexandros Stamatakis, Nick Goldman & Christophe Dessimoz

Final thoughts

Maximum likelihood and parsimony methods allow us to efficiently explore much larger datasets than currently possible with Bayesian methods.

More efficient tree building, dating, ancestral state reconstruction and subsampling strategies can be combined into efficient workflows for working with massive datasets.

Regardless of how big your dataset is, the methods we discussed allow for easier and more efficient data exploration strategies that can be run upstream of more computationally intensive analyses.