

Predicting the (very near) future: forecasting pathogen evolution

Molecular Epidemiology of Infectious Diseases

Lecture 14

April 20nd, 2026

“No scientific theory is worth anything unless it enables us to predict something which is actually going on. Until that is done, theories are a mere game of words, and not such a good game as poetry”

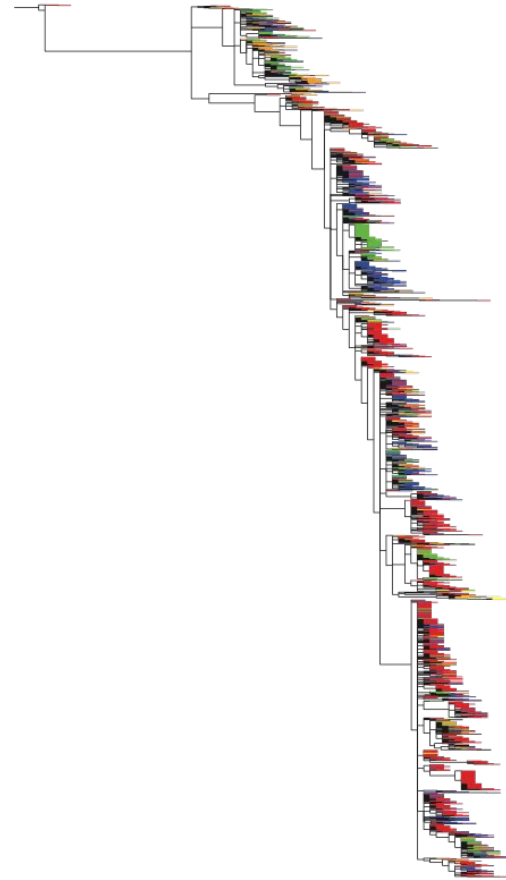
J.B.S Haldane (Adventures of a Biologist, 1937)

**Most of the
approaches we've
considered are
retrospective... can
we say anything
about the future?**

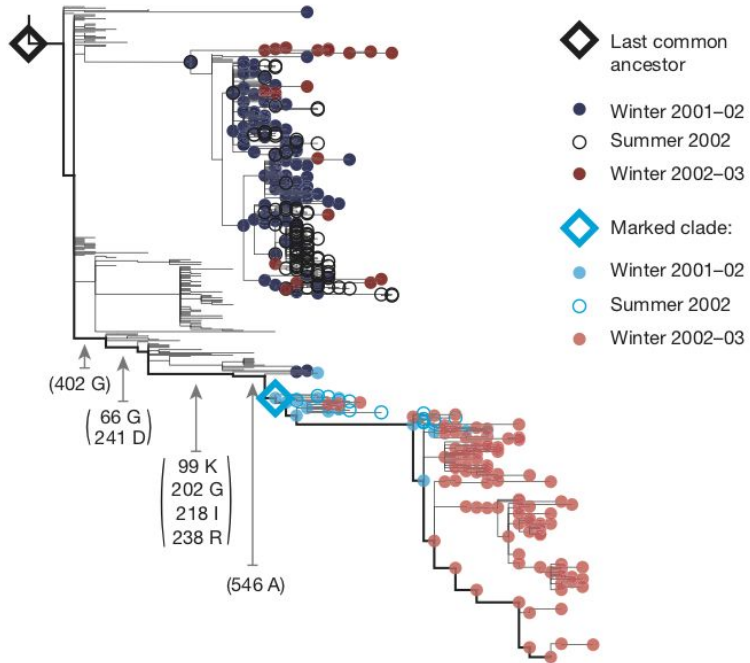
Influenza A (H3N2)

New antigenic variants periodically replace older strains:

- New antigenic variants emerge and escape antibody-based immunity against earlier strains.
- **Antigenic drift** leads to a ladder-like structure with a trunk lineage
- Flu vaccines need to be updated yearly to avoid antigenic mismatch.



Forecasting short-term flu evolution

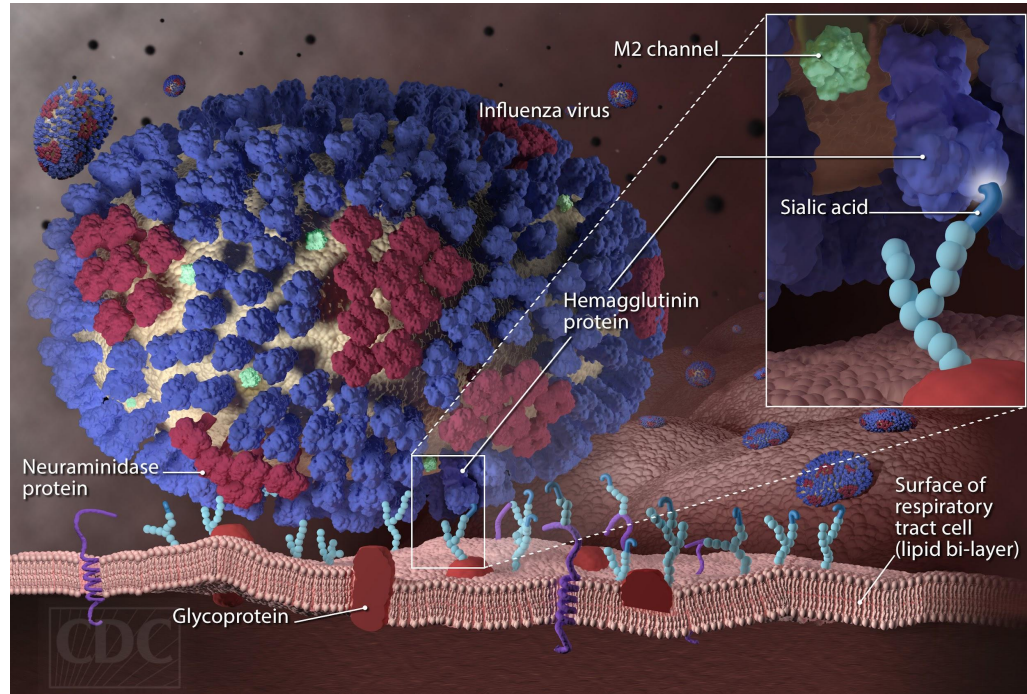


Consider the evolution dynamics of different influenza *clades*

The frequency X_v of a particular clade can be predicted based on the fitness f_i of individual strains i in a clade:

$$\hat{X}_v(t+1) = \sum_{i:v,t} x_i \exp(f_i)$$

Influenza hemagglutinin and cell entry



Forecasting short-term flu evolution

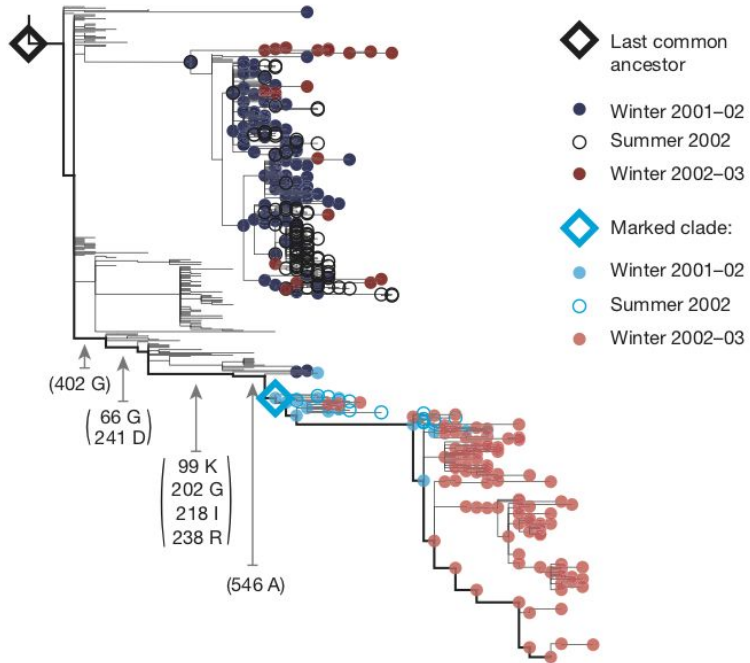
Luskza & Lassig (2014) consider two main factors that influence the fitness f_i of a strain:

- 1) The amplitude of cross-immunity $\mathbf{C}(\mathbf{a}_i, \mathbf{a}_j)$ between strain i and all other strains j that have previously circulated in the host population
- 2) The fitness cost $\mathbf{L}(\mathbf{a}_i)$ of deleterious mutations at non-antigenic sites

Their overall fitness mapping function is:

$$f_i = f_0 - \mathcal{L}(\mathbf{a}_i) - \sum_{j: t_j < t_i} x_j \mathcal{C}(\mathbf{a}_i, \mathbf{a}_j)$$

Forecasting short-term flu evolution

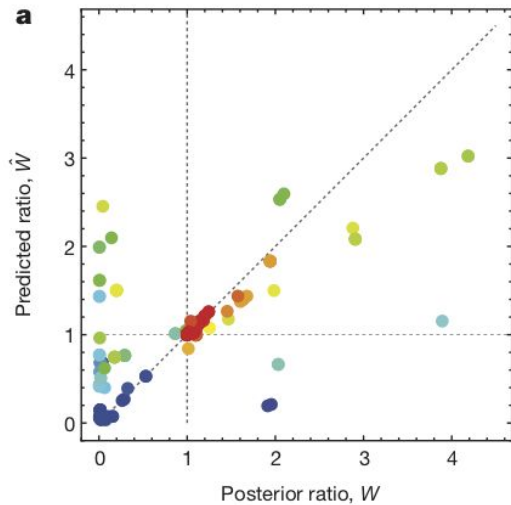


Consider the evolution dynamics of different influenza *clades*

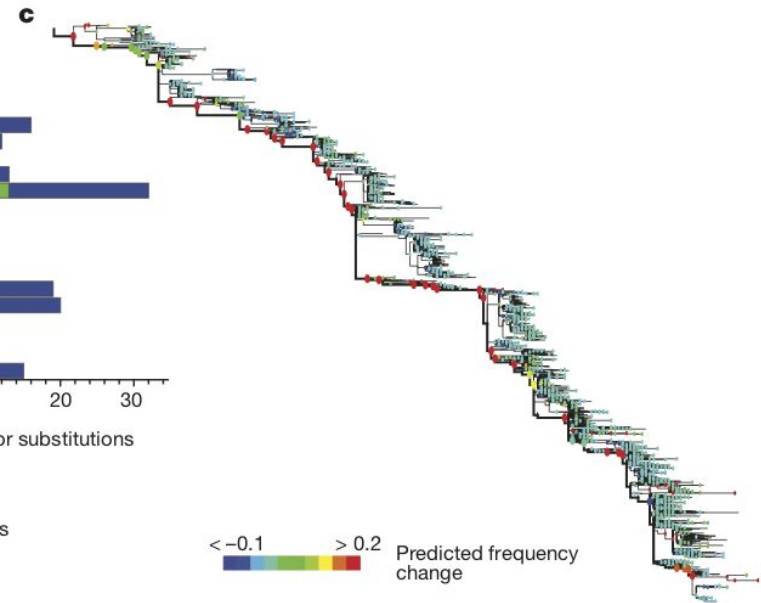
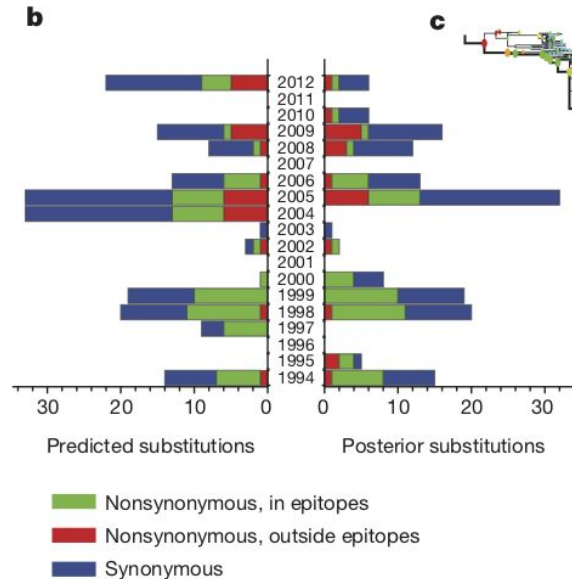
The frequency X_v of a particular clade can be predicted based on the fitness f_i of individual strains i in a clade:

$$\hat{X}_v(t+1) = \sum_{i:v,t} x_i \exp(f_i)$$

Forecasting short-term flu evolution

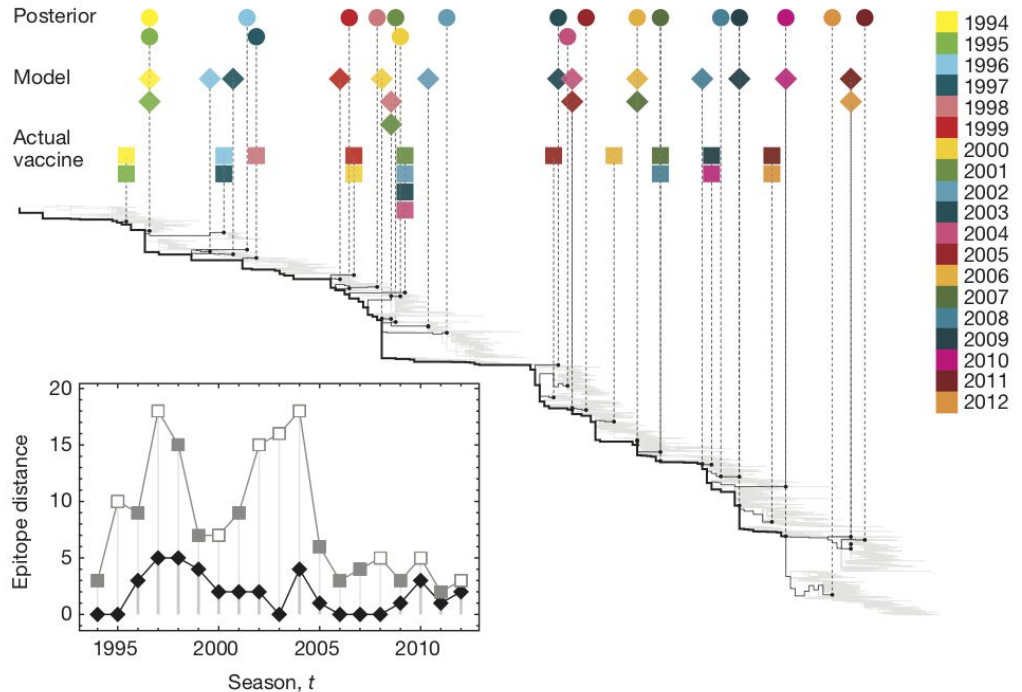


$$W_v = \frac{X_v(t+1)}{X_v(t)}$$



Forecasting short-term flu evolution

Evolutionary predictions can aid design of vaccines with optimal immunity to dominant strains in the next flu season.



**Can we predict
pathogen evolution
more generally?**

What do we need to know?

What mutations/genotypes are available?

Will the fate of new variants be determined by selection or drift?

How do genotypes map to fitness-related phenotypes?

How does fitness translate to epidemic potential at the population level?

What do we need to know?

What mutations/genotypes are available?

Will the fate of new variants be determined by selection or drift?

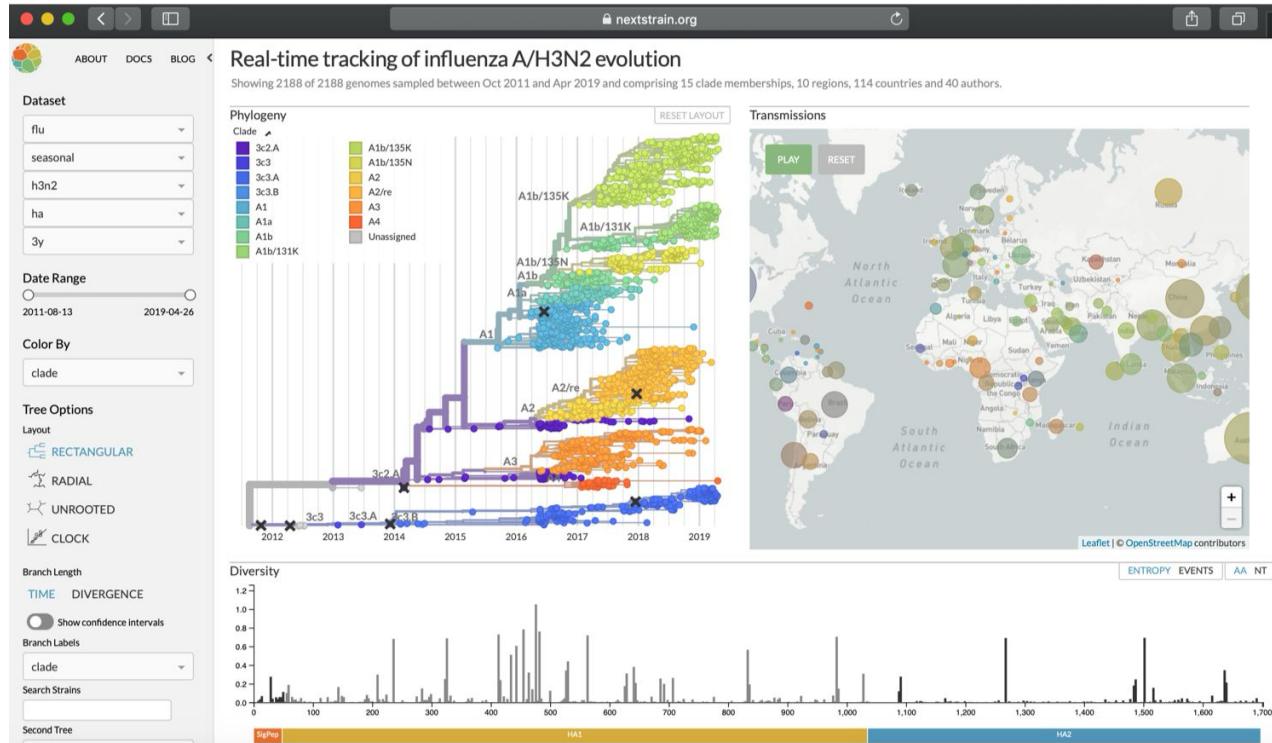
How do genotypes map to fitness-related phenotypes?

How does fitness translate to epidemic potential at the population level?

Mutational limits on prediction

At the very least, we need to know what mutations/genotypes are in a population to be able to predict anything about evolution.

Genomic surveillance



Mutational limits on prediction

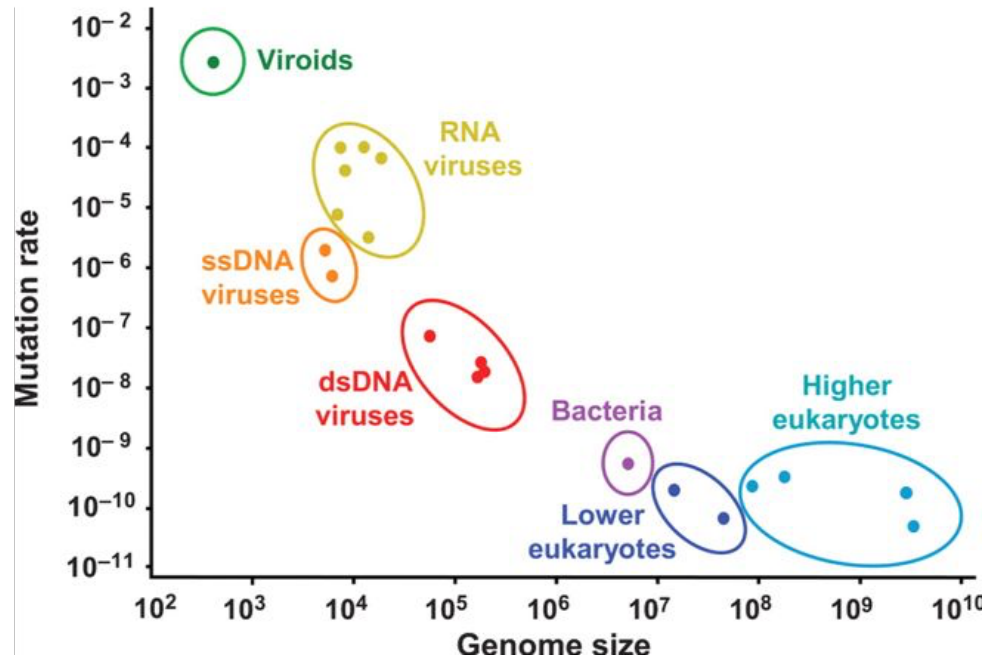
At the very least, we need to know what mutations/genotypes are in a population to be able to predict anything about evolution.

Meaningful predictions are probably limited to short-term predictions about standing genetic variation (or immediately accessible mutations).

Rapidly mutating microbes

Microbial evolution is often not mutation limited - high mutation rates and large population sizes ensure all possible mutations occur on short timescales.

Evolutionary predictions may then be extended to all locally accessible genotypes (e.g. genotypes one mutation away from existing strains).



Gago et al. (Science, 2009)

Mutational limits on prediction

At the very least, we need to know what mutations/genotypes are in a population to be able to predict anything about evolution

Meaningful predictions are probably limited to short-term predictions about standing genetic variation (or immediately accessible mutations).

For rapidly evolving microbial pathogens, we may be able to extend these predictions to all locally accessible genotypes.

Long-term predictions are limited by the stochastic nature of the mutation process and what mutations will enter a population

Mutational limits on prediction

At the very least, we need to know what mutations/genotypes are in a population to be able to predict anything about evolution

Meaningful predictions are probably limited to short-term predictions about standing genetic variation (or immediately accessible mutations).

For rapidly evolving microbial pathogens, we may be able to extend these predictions to all locally accessible genotypes.

Long-term predictions are however still limited by the stochastic nature of the mutational process and what combinations of mutations will occur in the future.

What do we need to know?

What mutations/genotypes are available?

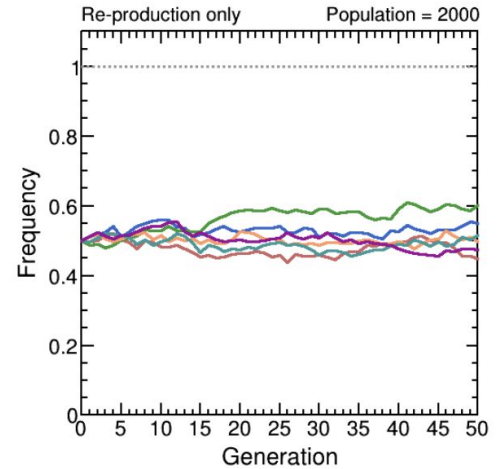
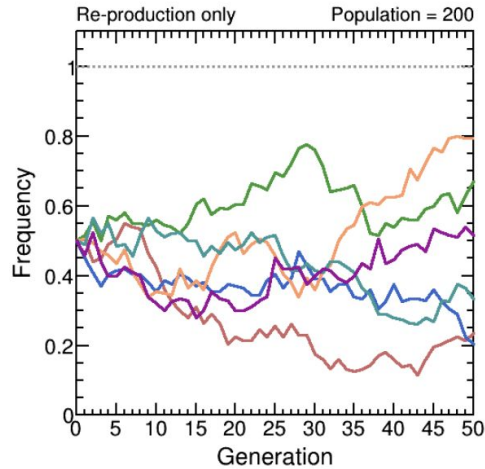
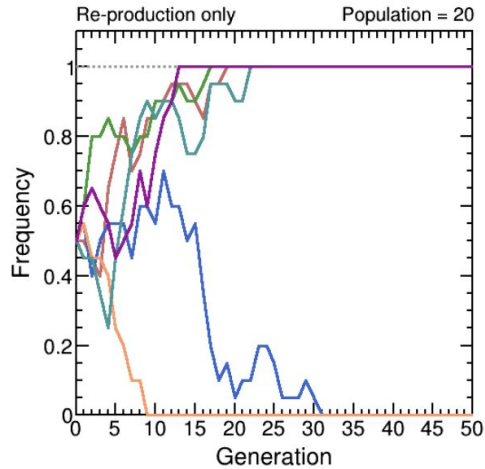
Will the fate of new variants be determined by selection or drift?

How do genotypes map to fitness-related phenotypes?

How does fitness translate to epidemic potential at the population level?

Genetic drift

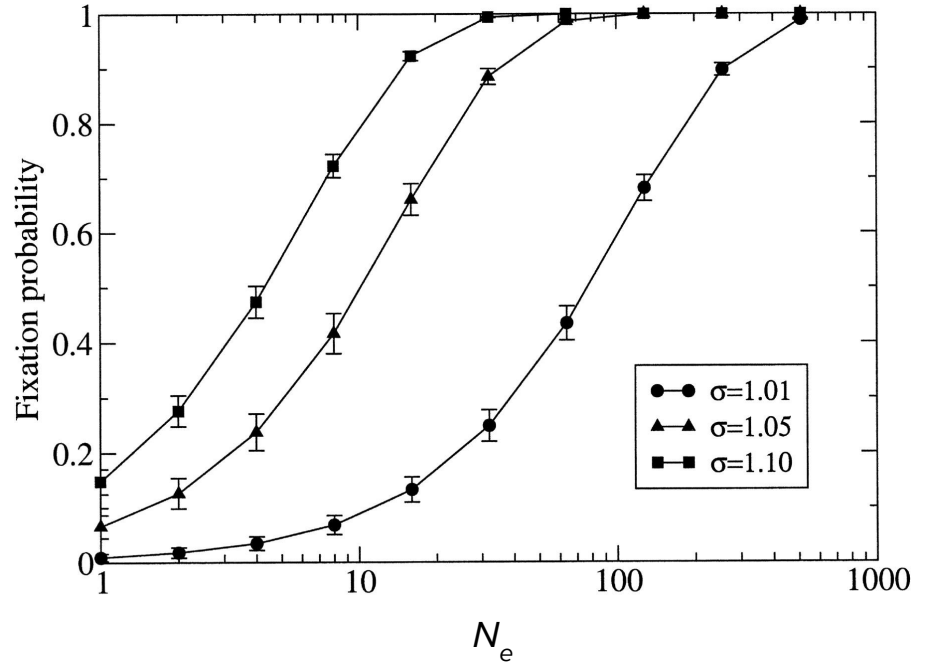
Genetic drift refers to stochastic fluctuations in genotype frequencies caused by random variation in reproduction and survival. Stochastic variation and drift play a larger role in smaller populations.



Genetic drift

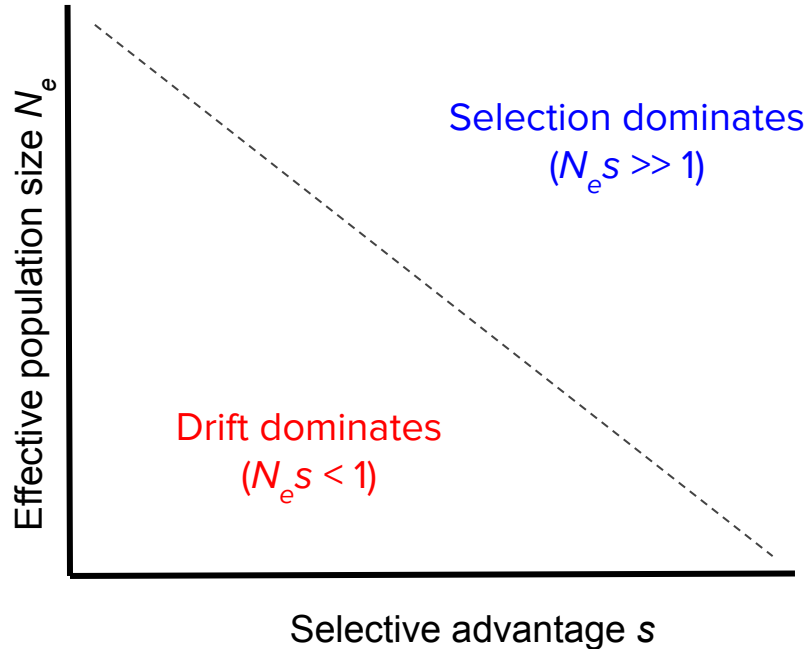
The probability that a beneficial mutation reaches fixation (freq \rightarrow 1.0) depends both on its selective advantage (s or σ) and the effective population size (N_e) – the number of individuals that contribute progeny to the next generation.

$$S = W_{mut} - W_{wt}$$



Selection vs. drift

The relative importance of selection versus drift is determined by $N_e s$. At low values of $N_e s$ drift will dominate selection, making prediction very difficult.

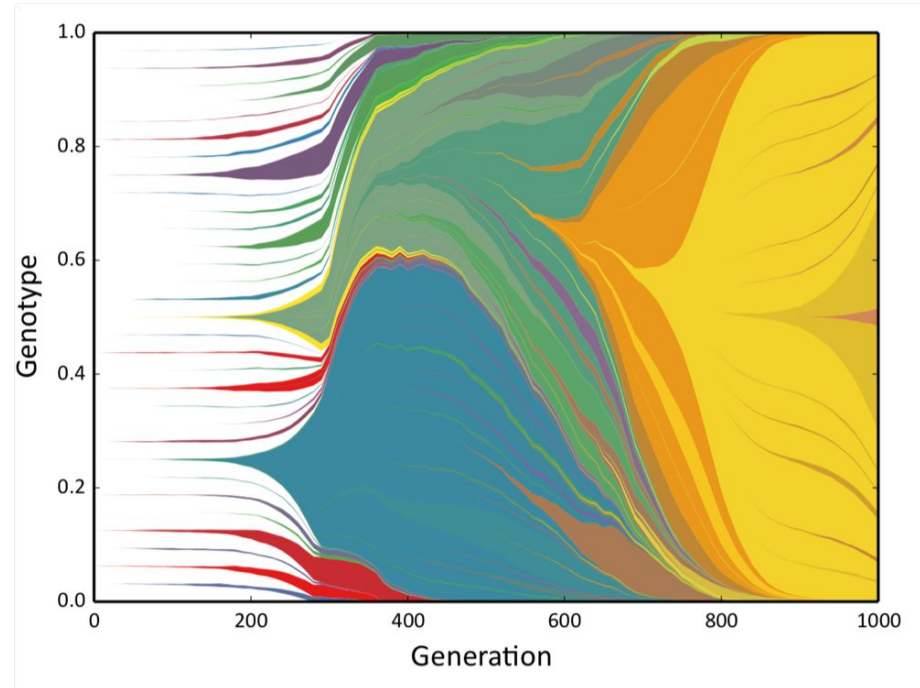


Clonal interference

Clonal interference arises in large asexual populations with high mutations rates.

Multiple lineages with beneficial mutations compete with one another.

Increases the odds that the most fit genotype goes to fixation.



Selection versus drift in pathogen pops

In large populations, the role of genetic drift becomes negligible relative to selection.

However, its less clear whether pathogen effective population sizes are large enough that we can ignore drift.

Clonal interference can enhance overall predictability: it increases odds of the most fit genotype going to fixation even if multiple mutations are required.

Evolution in large microbial populations may be more predictable than others!

What do we need to know?

What mutations/genotypes are available?

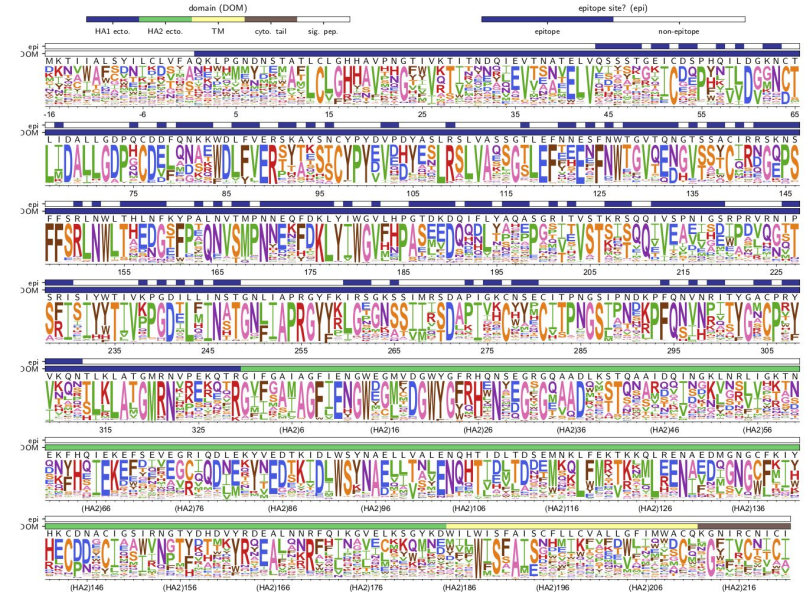
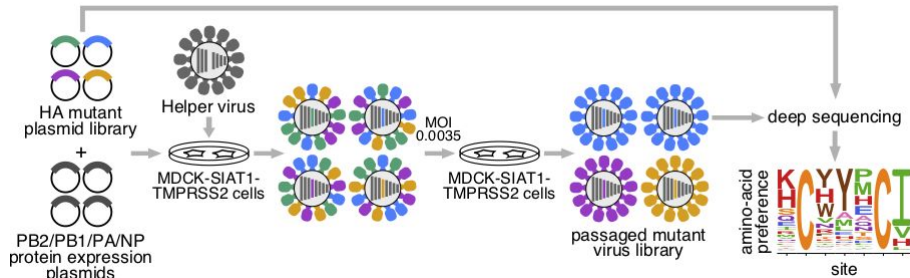
Will the fate of new variants be determined by selection or drift?

How do genotypes map to fitness-related phenotypes?

How does fitness translate to epidemic potential at the population level?

Deep mutational scanning

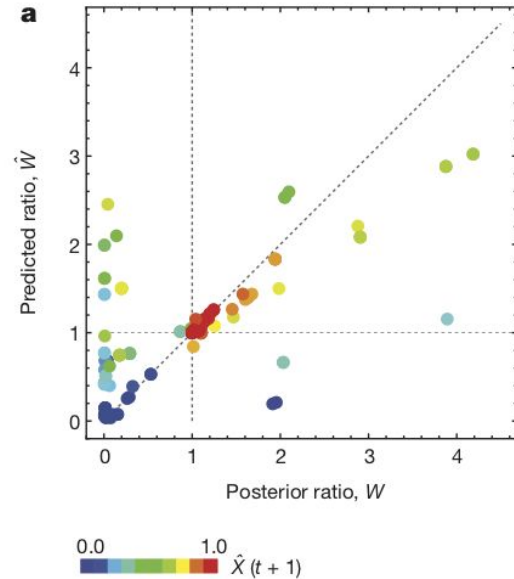
Reverse genetics approaches can be used to systematically explore the genotype to phenotype map using large libraries of mutants.



But genetic context matters too

Luskza and Lassig found the models that only consider “adaptive” changes in epitope regions are 40% less accurate than models that all consider changes in background fitness due to deleterious mutations in other parts of the genome.

$$f_i = f_0 - \mathcal{L}(\mathbf{a}_i) - \sum_{j: t_j < t_i} x_j \mathcal{C}(\mathbf{a}_i, \mathbf{a}_j)$$



Context dependence

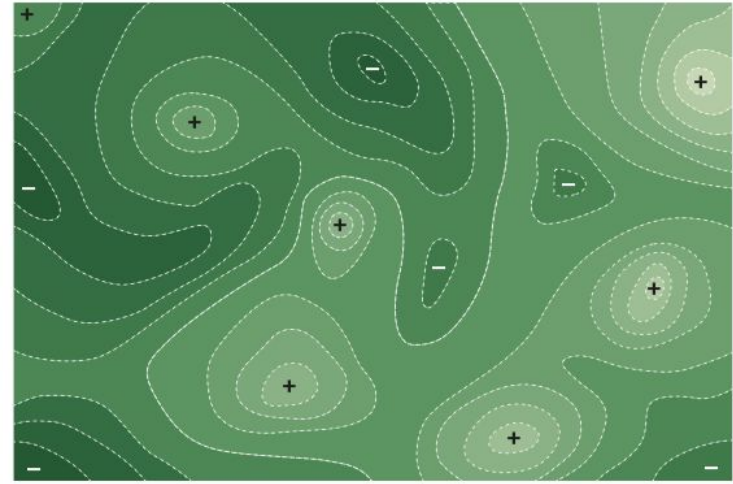
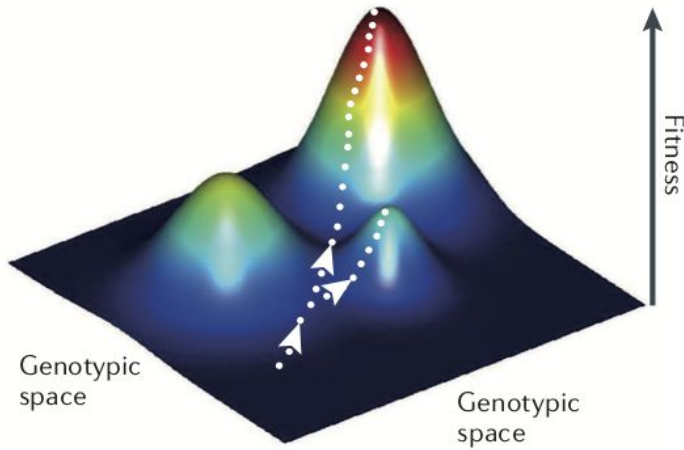
How predictable phenotypes/fitness are based on genotypes largely depends on whether phenotypes are context dependent:

Epistasis: dependence on genetic background including interactions among mutations.

Pleiotropy: the effects of mutations on multiple traits or the same trait across different environments.

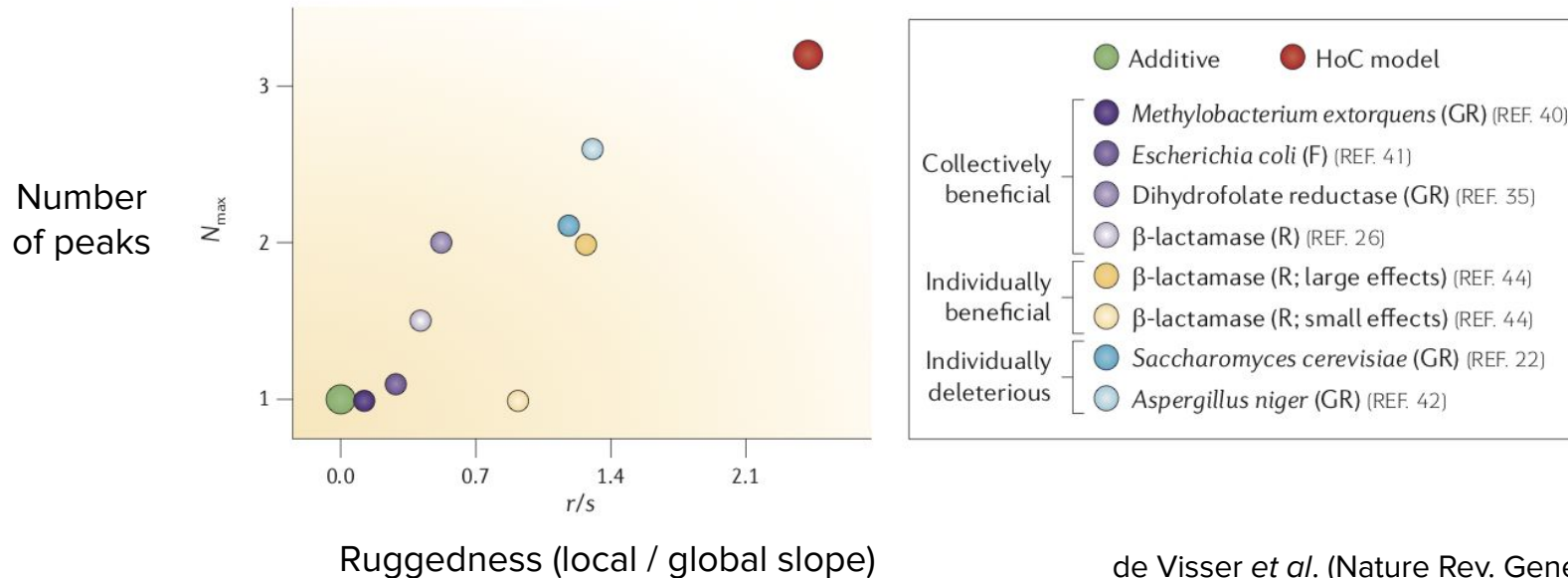
Epistasis in fitness landscapes

Epistasis largely controls the smoothness/ruggedness of the fitness landscape. Strong epistasis makes prediction difficult due to rugged landscapes.



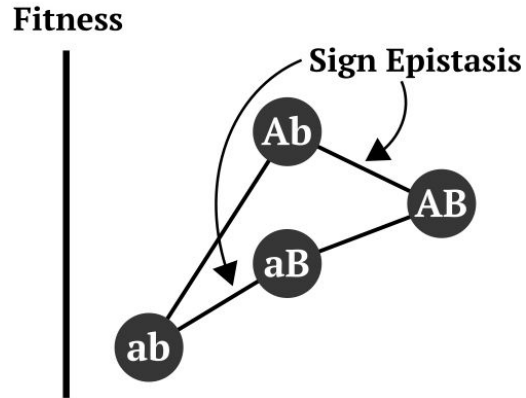
Epistasis in fitness landscapes

Empirical fitness landscapes tend to have intermediate levels of ruggedness.

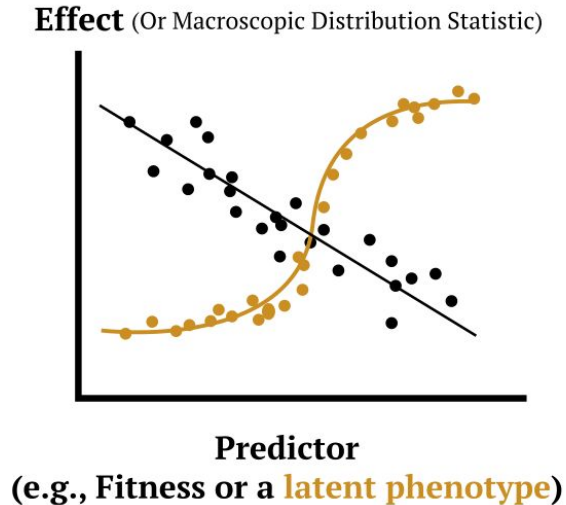


Two types of epistasis

Idiosyncratic Epistasis

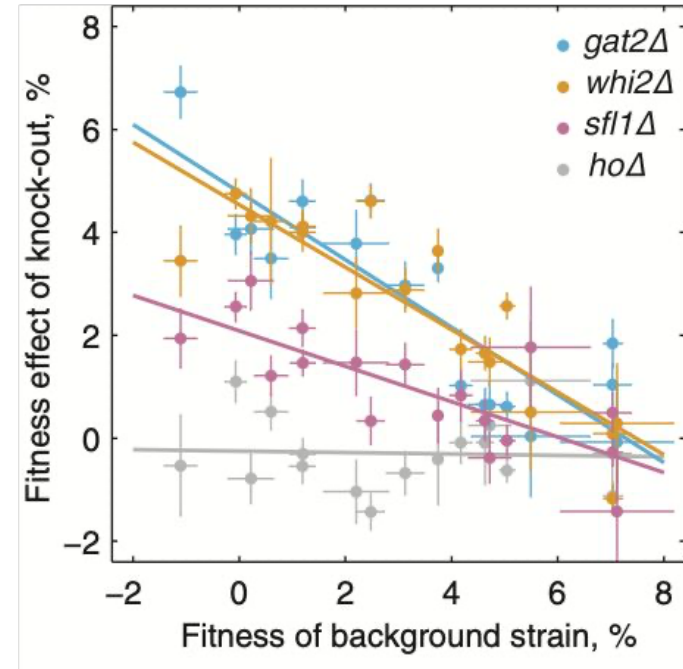


Global Epistasis



Global epistasis

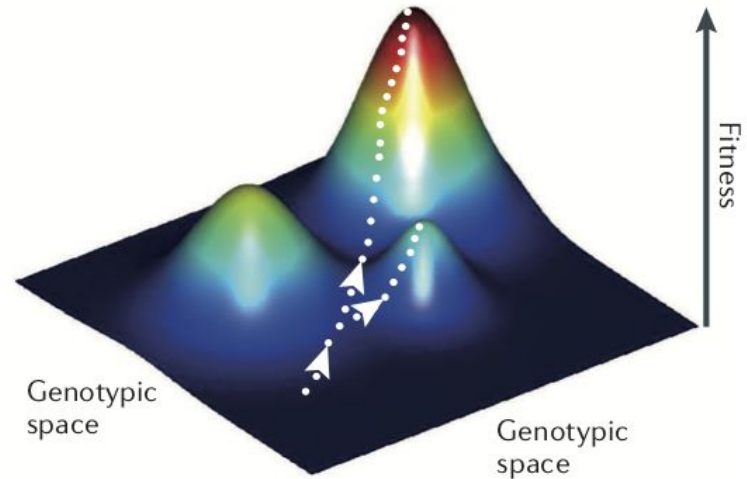
Mutations often exhibit *global epistasis* where their fitness effects depend on starting fitness but are “independent of the specific identify of mutations present in the background”.



Global epistasis

Mutations often exhibit *global epistasis* where their fitness effects depend on starting fitness but are “independent of the specific identify of mutations present in the background”.

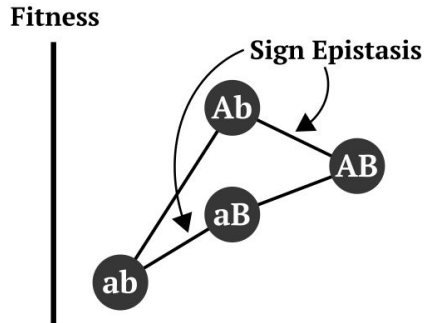
This is often seen as “diminishing returns” on the effects of beneficial mutations in already fit genotypes.



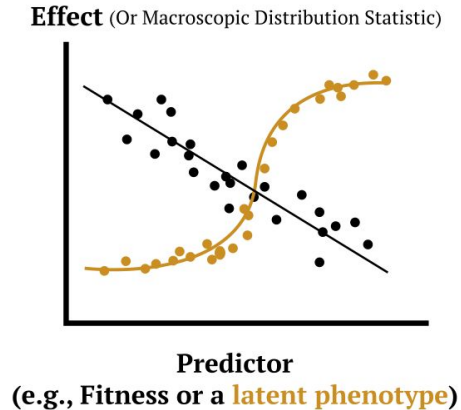
Two types of epistasis

Overall, idiosyncratic epistasis makes the fitness effects of mutations less predictable whereas global epistasis makes fitness effects more predictable.

Idiosyncratic Epistasis



Global Epistasis



What do we need to know?

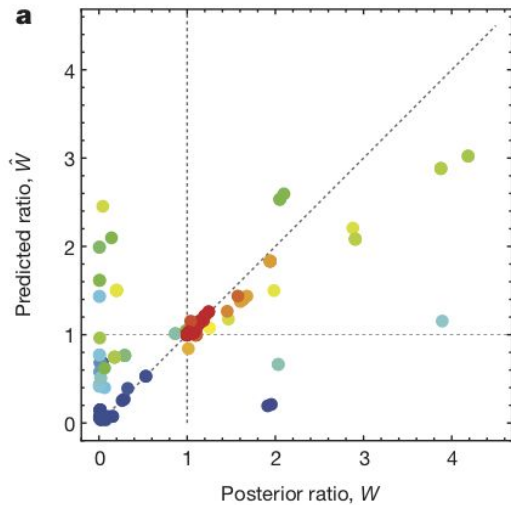
What mutations/genotypes are available?

Will the fate of new variants be determined by selection or drift?

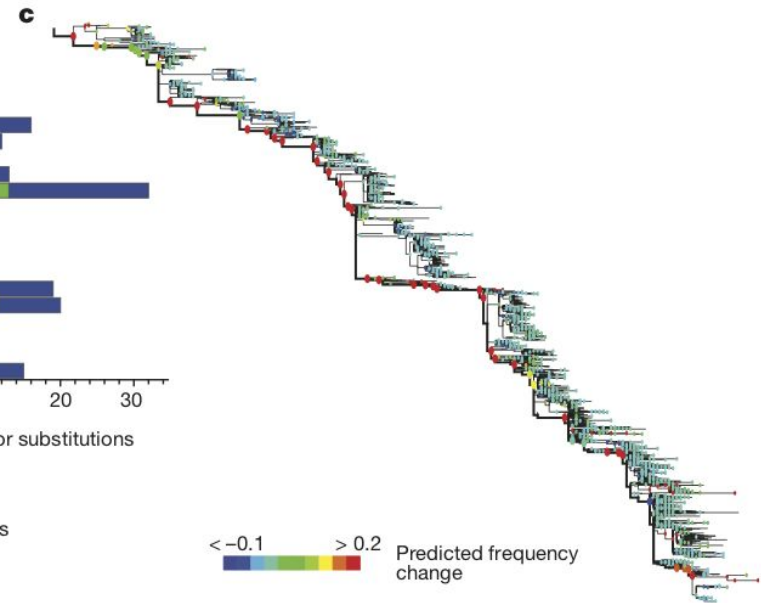
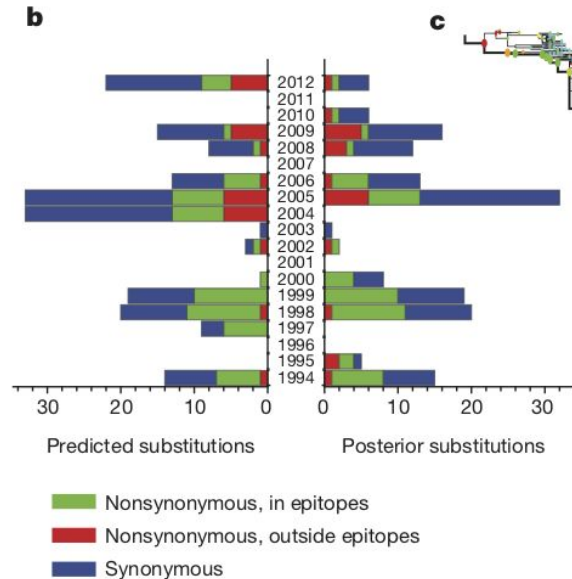
How do genotypes map to fitness-related phenotypes?

How does fitness translate to epidemic potential at the population level?

Forecasting short-term flu evolution



$$W_v = \frac{X_v(t+1)}{X_v(t)}$$

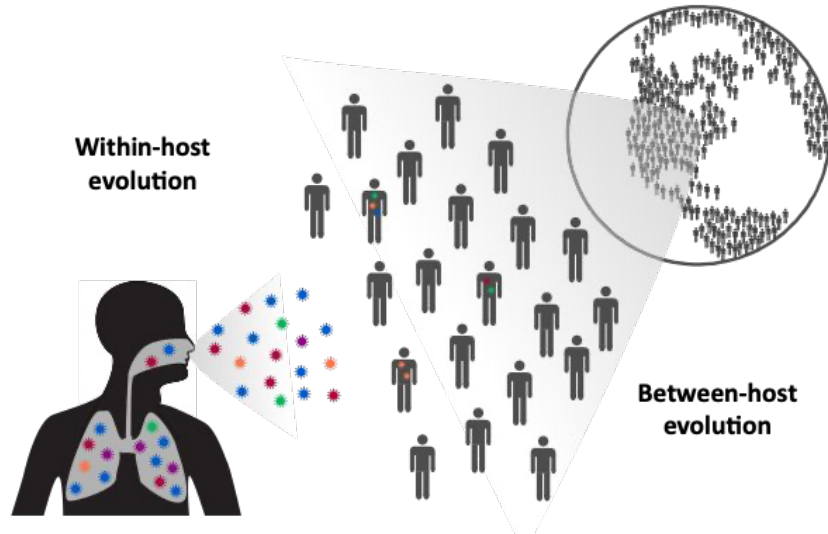


“Any prediction of evolution is essentially an estimate of fitness differences between strains”

Luksza & Lassig (2014)

Translating between scales

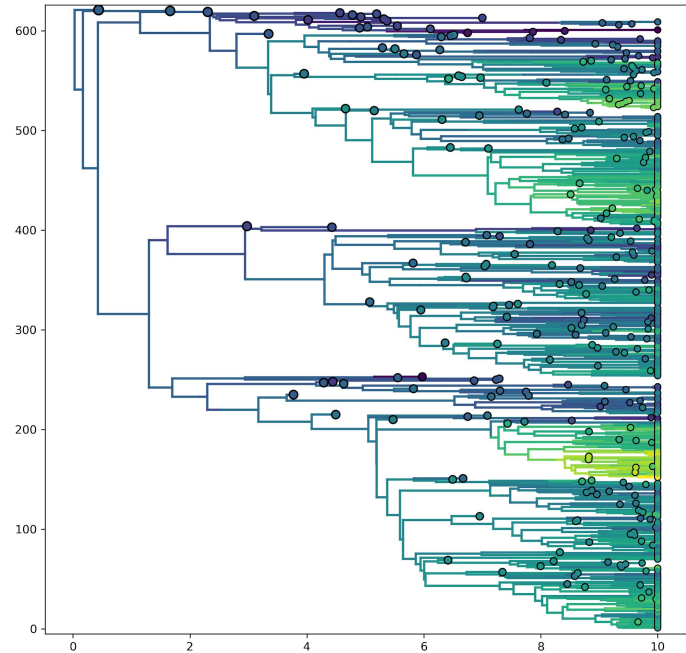
To make accurate predictions we need to know how pathogen phenotypes related to within-host fitness translate to population-level fitness between hosts.



Fitness shapes pathogen phylogenies

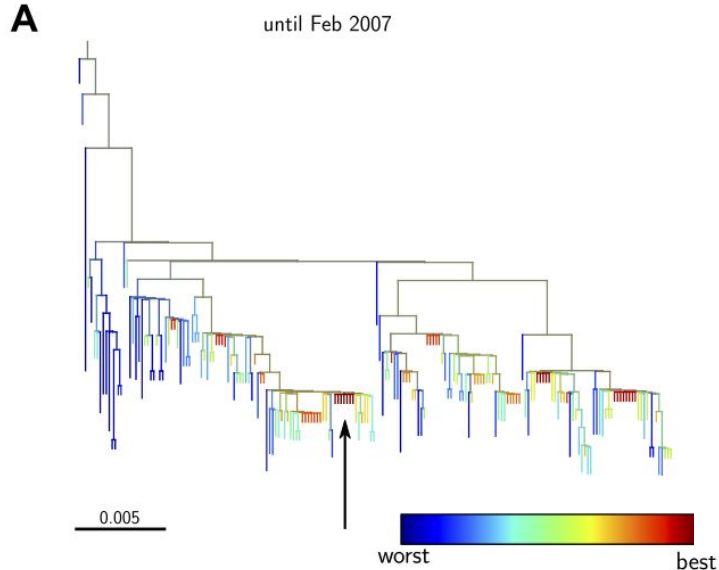
More fit lineages will have higher growth rates and therefore branch more often... leaving behind more sampled descendents in a phylogeny.

branching = birth/transmission events



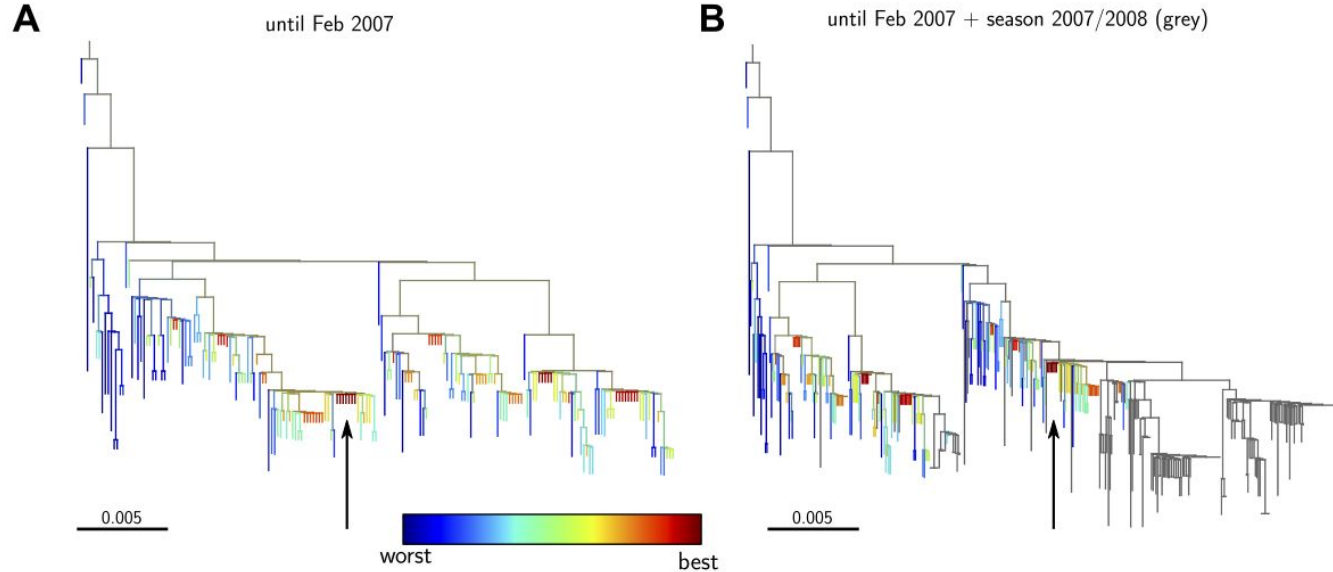
Predicting evolution from tree shape

Branching rates in pathogen phylogenies correlate strongly with fitness



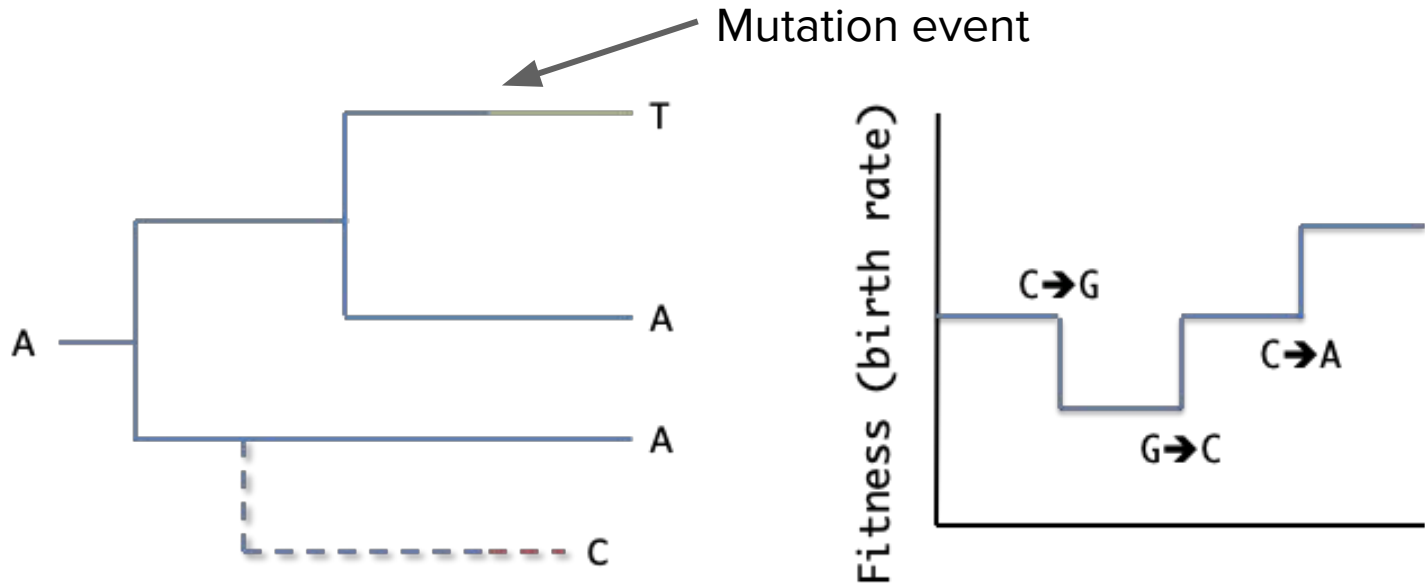
Predicting evolution from tree shape

Branching rates in pathogen phylogenies correlate strongly with fitness



Multi-type birth-death models

Allows for different types of individuals (e.g. genotypes) that can vary in their birth or death rates and therefore their fitness values.



Fitness of HIV drug resistance mutations

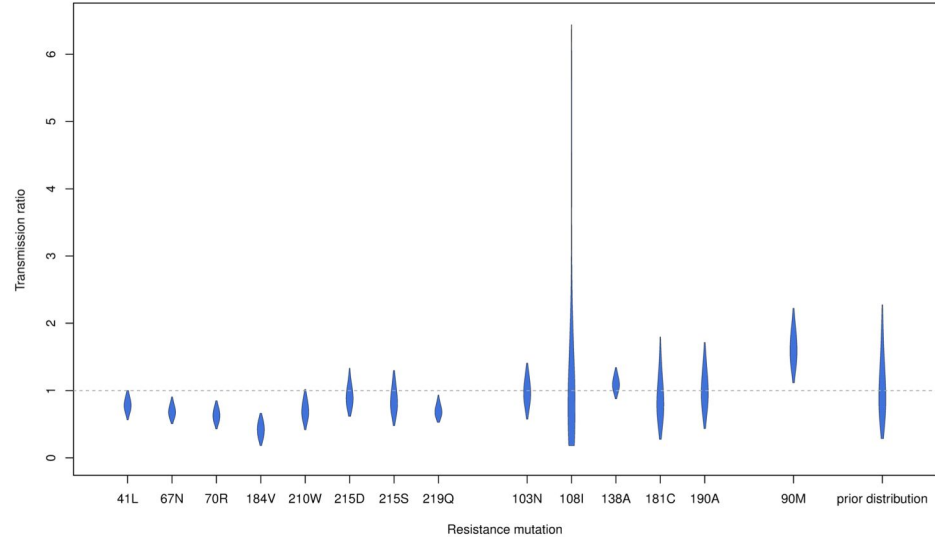
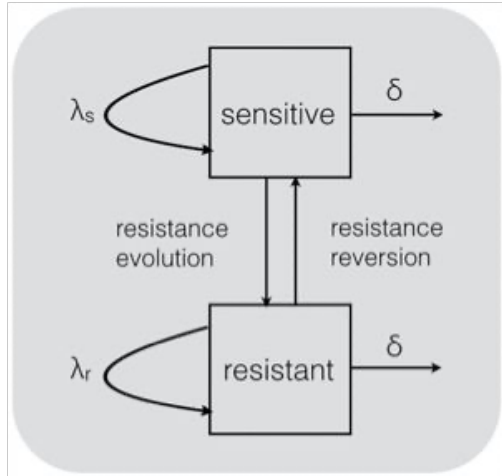


Table 1. Resistance mutations with numbers of corresponding clusters and samples, related drugs and drug usage dates within Switzerland.

Resistance mutation	nRTI								NNRTI				PI		
	41L	67N	70R	184V	210W	215D	215S	215Y	219Q	103N	108I	138A	181C	190A	90M
Number (#) of clusters of size ≥ 2	56	23	19	35	18	18	16	25	20	25	10	46	8	8	14
# Sequences in clusters	927	667	712	1011	481	569	494	807	605	725	334	1014	329	311	389
# Resistant samples in clusters	93	39	26	44	26	41	31	28	28	38	11	109	10	12	38
Drug (SHCS drug codes)	AZT D4T	AZT D4T	AZT D4T	3TC ABC FTC	AZT D4T	AZT D4T	AZT D4T	AZT D4T	AZT D4T	NVP EFV	NVP EFV	RPV	NVP EFV ETV RPV	NVP EFV	NVP SQV
Drug usage $\geq 1\%$	1987	1987	1987	1995.5	1987	1987	1987	1987	1987	1997	1997	2013	1997	1997	1996
Drug usage $< 1\%$	-	-	-	-	-	-	-	-	-	-	-	-	-	-	2008

A pathogen's fitness is a composite phenotype determined by many different intrinsic and extrinsic factors.

Forecasting SARS-CoV-2 evolution

Cell

 CellPress
OPEN ACCESS

Article

Population immunity predicts evolutionary trajectories of SARS-CoV-2

Matthijs Meijers,¹ Denis Ruchnewitz,¹ Jan Eberhardt,¹ Marta Łuksza,² and Michael Lässig^{1,3,*}

¹Institute for Biological Physics, University of Cologne, Zùlpicherstr. 77, 50937 Kùln, Germany

²Tisch Cancer Institute, Departments of Oncological Sciences and Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA

³Lead contact

*Correspondence: mlaessig@uni-koeln.de

<https://doi.org/10.1016/j.cell.2023.09.022>

Forecasting SARS-CoV-2 evolution

Cell

CellPress
OPEN ACCESS

Article

Population immunity predicts evolutionary trajectories of SARS-CoV-2

Matthijs Meijers,¹ Denis Ruchnewitz,¹ Jan Eberhardt,¹ Marta Łuksza,² and Michael Lässig^{1,3,*}

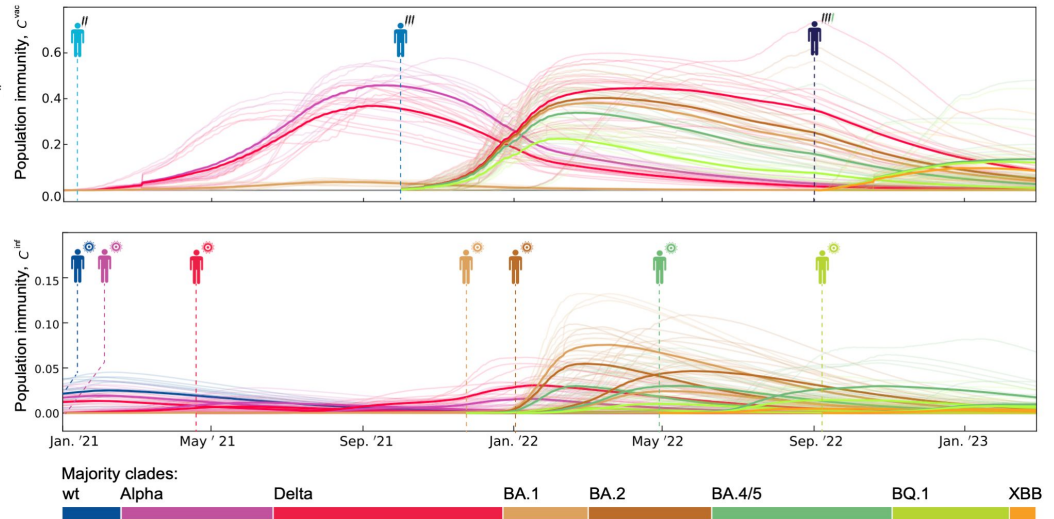
¹Institute for Biological Physics, University of Cologne, Zùlpicherstr. 77, 50937 Kùln, Germany

²Tisch Cancer Institute, Departments of Oncological Sciences and Genetics and Genomic Sciences, Icahn School of New York, NY, USA

³Lead contact

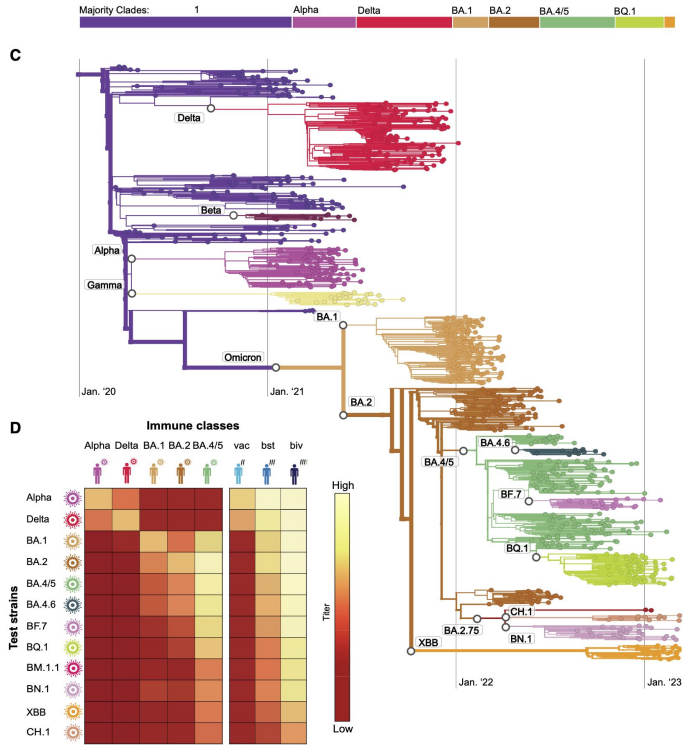
*Correspondence: mlassig@uni-koeln.de

<https://doi.org/10.1016/j.cell.2023.09.022>



Meijers *et al.*, (Cell, 2023)

Forecasting SARS-CoV-2 evolution



Meijers et al. use a fitness prediction model very similar to Luskza & Lassig:

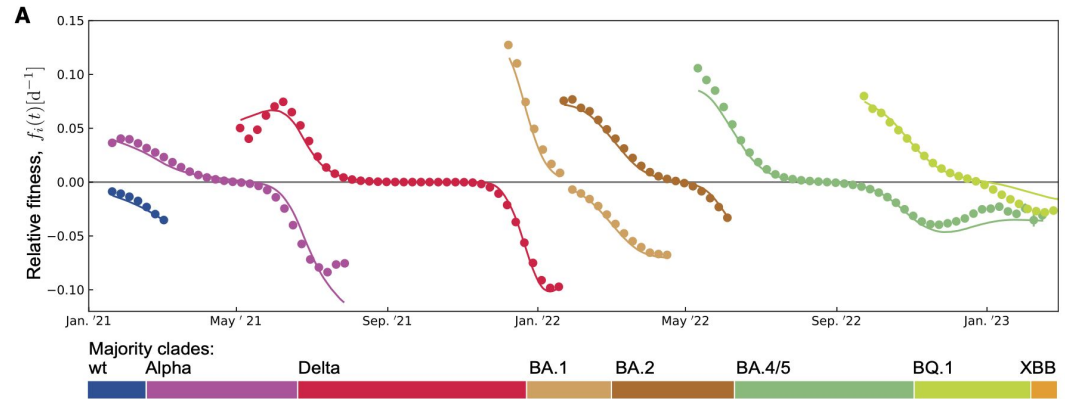
$$f_i = f_0 - \mathcal{L}(\mathbf{a}_i) - \sum_{j: t_j < t_i} x_j \mathcal{C}(\mathbf{a}_i, \mathbf{a}_j)$$

Fitness estimated from growth rates of individual variants allows them to forecast near-term changes in variant frequencies:

$$\hat{X}_v(t+1) = \sum_{i:v,t} x_i \exp(f_i)$$

Forecasting SARS-CoV-2 evolution

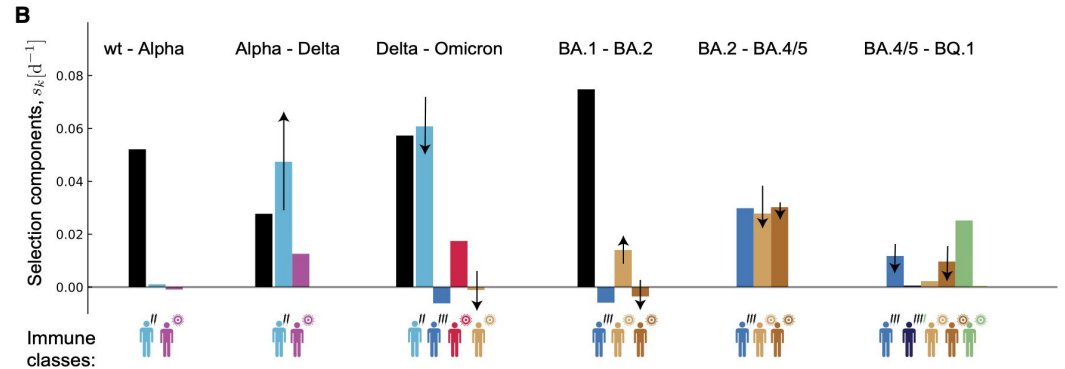
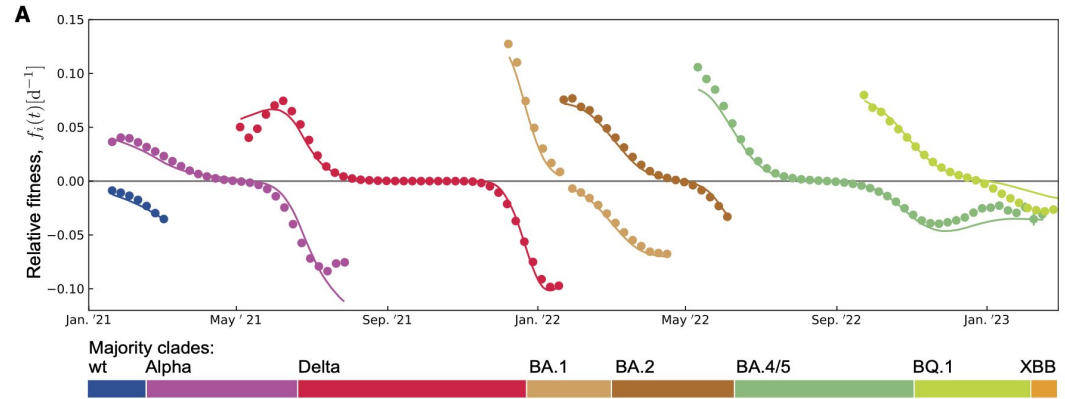
Fitting the model allows them to estimate the time-varying fitness of each variant as a function of other variant's current and past prevalence...



Forecasting SARS-CoV-2 evolution

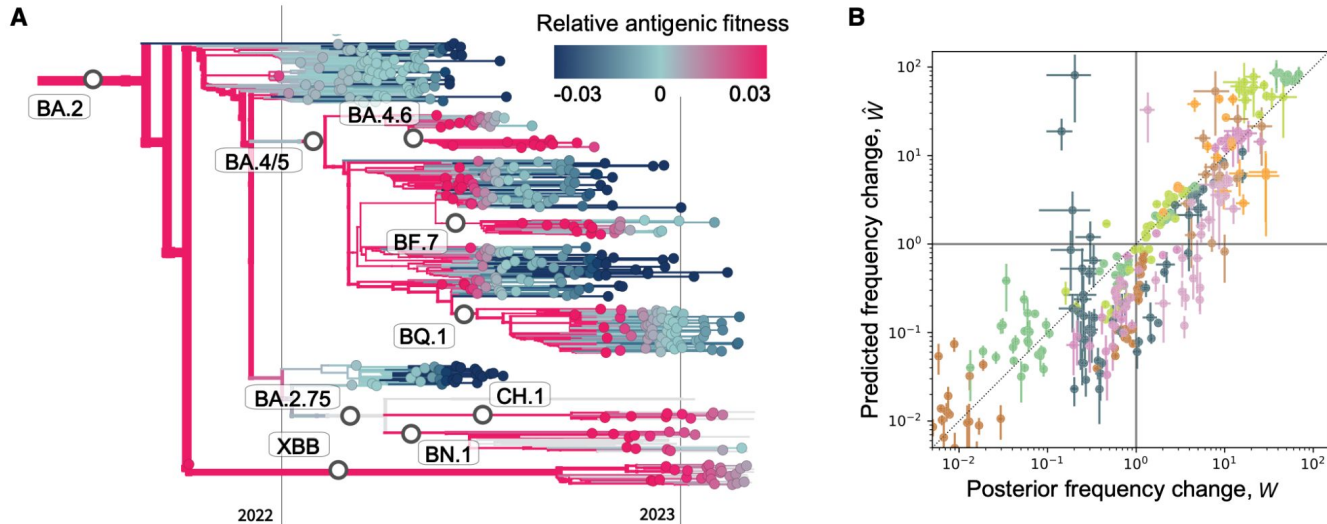
Fitting the model allows them to estimate the time-varying fitness of each variant as a function of other variant's current and past prevalence...

As well as compute the strength of selection (s) acting on different components of viral fitness.



Forecasting SARS-CoV-2 evolution

Estimating the relative fitness of competing variants in terms of both intrinsic and antigenic fitness allows for variant frequencies to be predicted quite accurately.



What do we need to know?

What mutations/genotypes are available?

Will the fate of new variants be determined by selection or drift?

How do genotypes map to fitness-related phenotypes?

How does fitness translate to epidemic potential at the population level?

**Can we predict
pathogen evolution
more generally?**

Analogy: Forecasting the weather

Despite the fact that the physical models required to predict the weather were developed in the 19th century, it still took another hundred years for reliable forecasts to emerge because of the need for massive amounts of atmospheric data and computing power.

But once short-term forecasts could be made, methods could be iteratively tested and improved, and forecasting advanced remarkably quickly.

A brief history of weather forecasting:

<https://www.newyorker.com/magazine/2019/07/01/why-weather-forecasting-keeps-getting-better>

The future of evolutionary predictions

We have the theory, methods and data to predict short-term evolution

- Predictive genotype-to-fitness models
- High-throughput phenotypic data
- Genomic surveillance data
- Predictive evolutionary/epidemiological models

We will likely get it wrong many times before we get it right but the fact that we can repeatedly test predictions on short timescales means that we can iteratively and rapidly improve our evolutionary forecasts.

In class discussion on Wednesday

Please read these two papers for class on Wednesday:

Wortel, M. T., Agashe, D., Bailey, S. F., Bank, C., Bisschop, K., Blankers, T., ... & Pennings, P. S. (2023). Towards evolutionary predictions: Current promises and challenges. *Evolutionary Applications*, 16(1), 3-21.

Meijers, M., Ruchnewitz, D., Eberhardt, J., Łuksza, M., & Lässig, M. (2023). Population immunity predicts evolutionary trajectories of SARS-CoV-2. *Cell*, 186(23), 5151-5164.

In class discussion on Wednesday

After you read these papers, please think about and be prepared to discuss:

1. How predictable is evolution in your favorite host-pathogen system?
2. What information is needed to make accurate predictions?
3. What is the time horizon of predictability?
4. What factors promote or limit predictability?
5. What is the biggest source of uncertainty surrounding predictions?