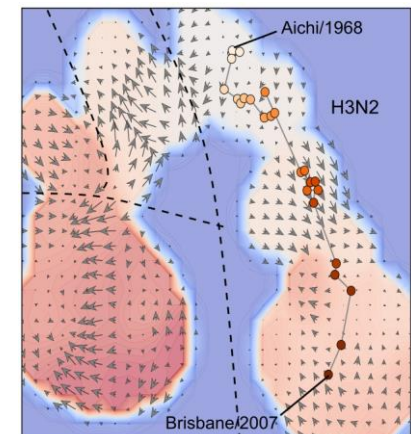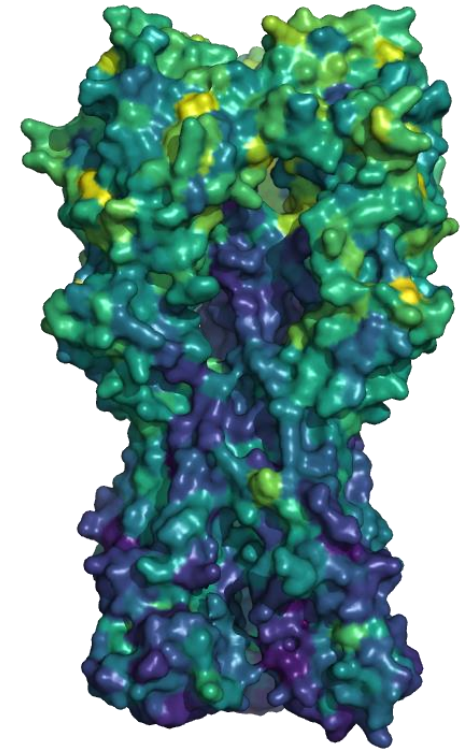# Predicting pathogen evolution with neural language models

Brian Hie
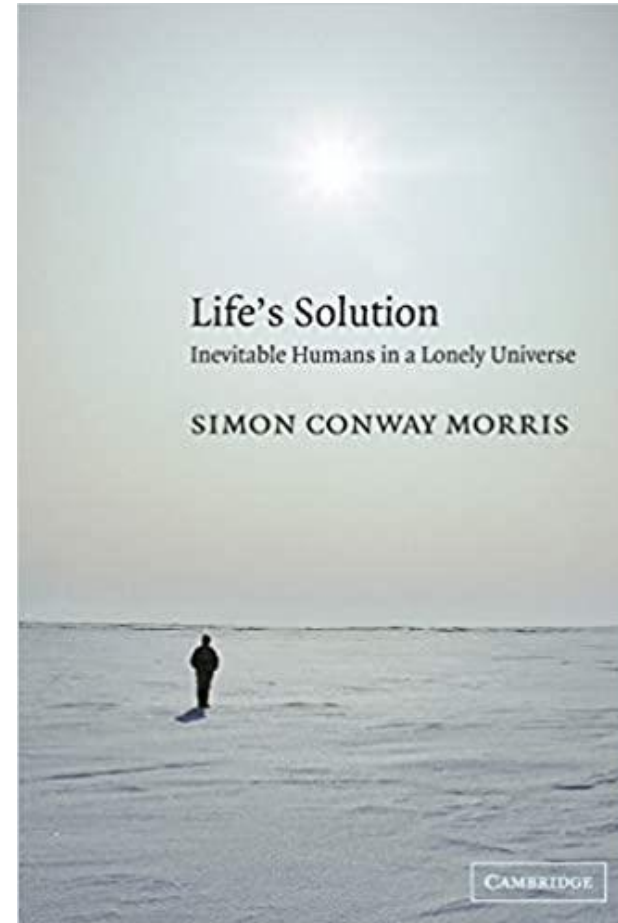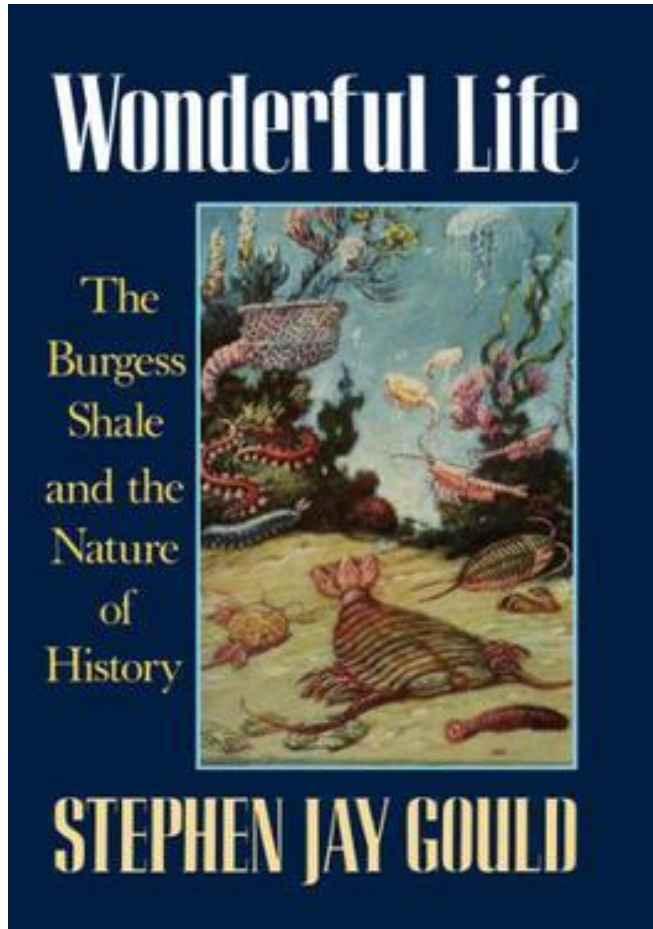
with Bonnie Berger, Bryan Bryson, Peter Kim,
Kevin Yang, and Ellen Zhong
October 28, 2021

# How predictable is evolution?

1. Learning the language of viral evolution and escape

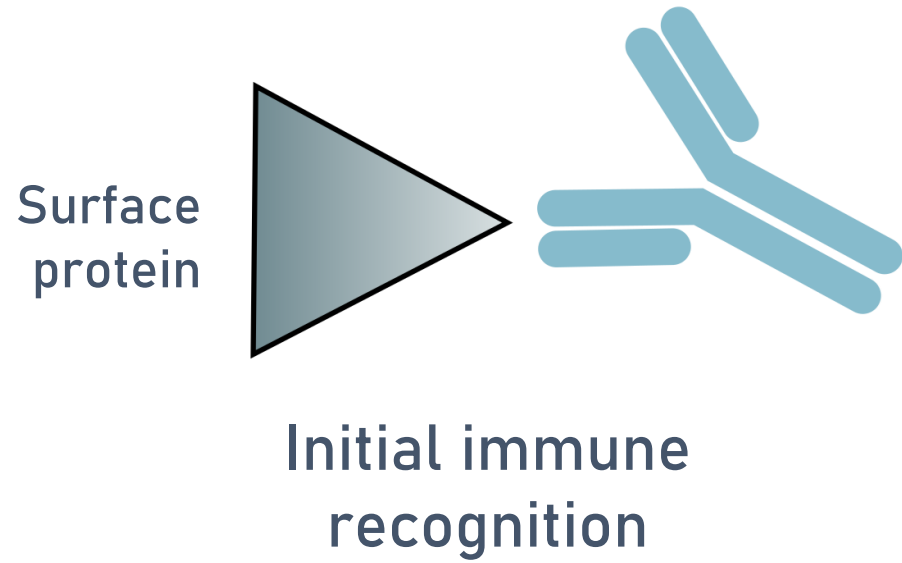2. Recovering evolutionary dynamics with "evolutionary velocity"

3. Looking forward

1. **Learning the language of viral evolution and escape**

2. Recovering evolutionary dynamics with "evolutionary velocity"
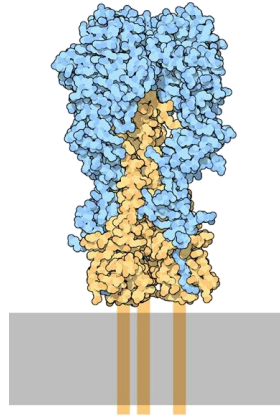
3. Looking forward

# Viral escape

Surface
protein

Initial immune
recognition

# Viral escape

Surface
protein

Initial immune
recognition

Protein mutation

Immune escape
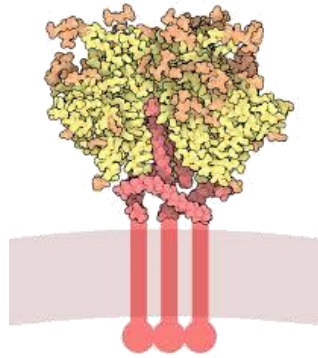
# Viral escape is a big problem

### Influenza A HA

### HIV Env

**Influenza**
- 250K–600K deaths a year
- Yearly vaccine that is 20–50% effective

**AIDS**
- 700K–1.2M deaths a year
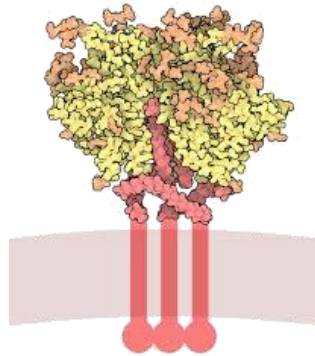- No effective vaccine

# Viral escape is a big problem

**Influenza A HA**

**HIV Env**

**SARS-CoV-2 Spike**
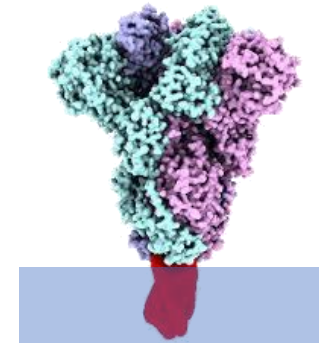
Influenza
- 250K–600K deaths a year
- Yearly vaccine that is 20–50% effective

AIDS
- 700K–1.2M deaths a year
- No effective vaccine

COVID-19
- 4.9M+ deaths
- Questions about durability of vaccine protection

# Small changes can have big semantic effects

# Small changes can have big semantic effects
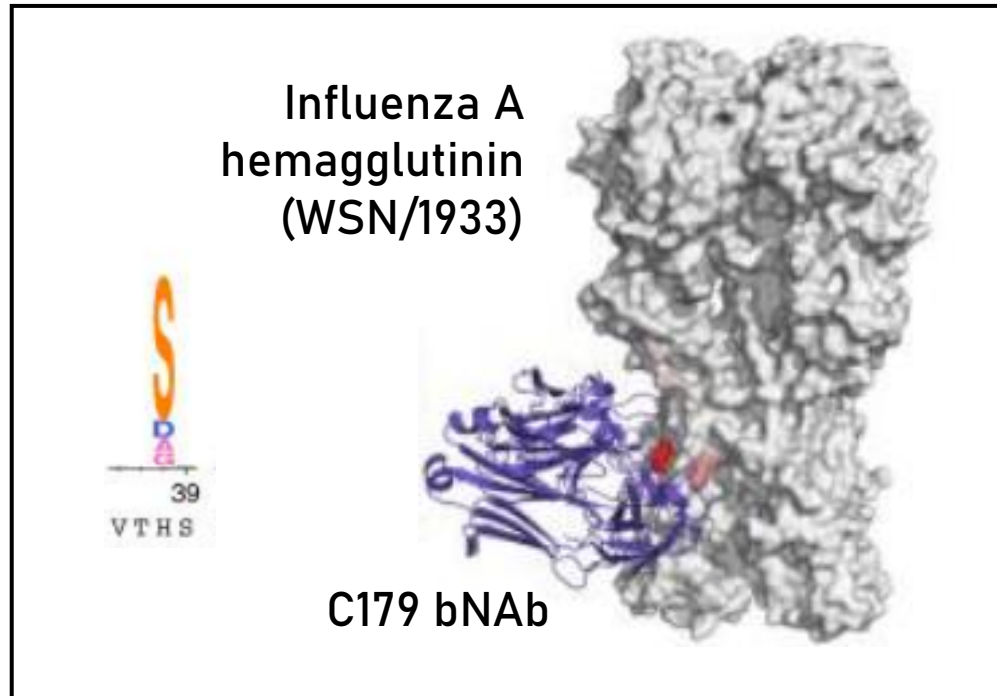
The boy pats the dog.

# Small changes can have big semantic effects

The boy pats the dog.

vs

The boy eats the dog.
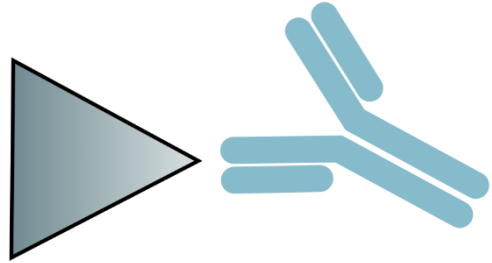
# Single residue change enables viral escape



Influenza A hemagglutinin (WSN/1933)

C179 bNAb

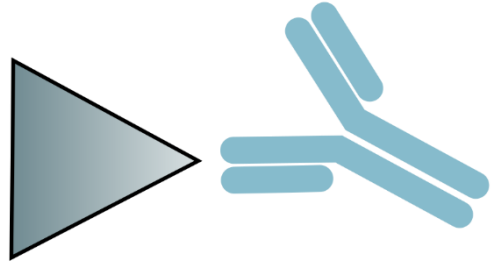From Doud, Lee, and Bloom. *Nat. Comm.* (2018)

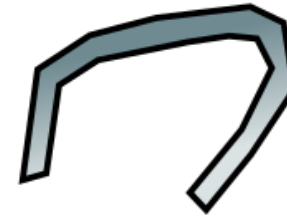H → S mutation means C179 no longer binds

# The language of viral escape

The boy pats the dog.

# The language of viral escape

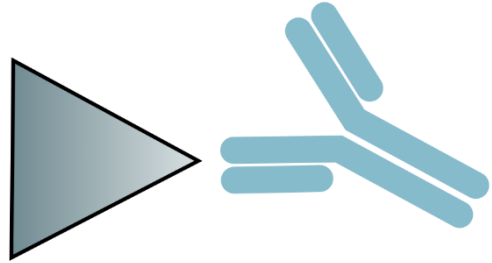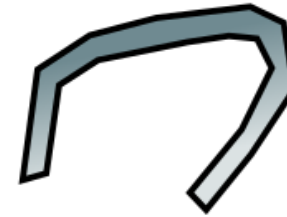

The boy pats the dog.

The boy patx the dog.

# The language of viral escape



The boy pats the dog.

The boy pat**x** the dog.

The boy p**e**ts the dog.

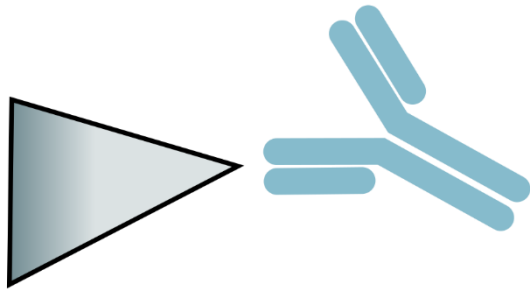# The language of viral escape
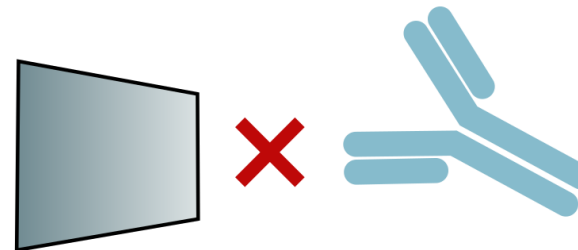


The boy pats the dog.

The boy pat**x** the dog.

The boy p**e**ts the dog.

The boy **e**ats the dog.

# Constrained semantic change search (CSCS)

- You're given a sequence of tokens from some language

- Goal: Find the single token change that:

1. Induces the largest semantic change

2. Is constrained by the rules/grammar of that language

# Constrained semantic change search (CSCS)

• You're given a sequence of tokens from some language

• Goal: Find the single token d

1. Induces the largest semar

2. Is constrained by the rule

For example:

A sequence of words from an English sentence

Or

A sequence of amino acids from a viral protein

# Some real example changes...

Original headline:

australian dead in bali

# Some real example changes...

Original headline:

australian dead in bali

Semantically closest:

<u>aussie</u> dead in bali

# Some real example changes...

Original headline:
australian dead in bali

Semantically closest:
<u>aussie</u> dead in bali

CSCS proposed change:
australian <u>ballet</u> in bali

# Some real example changes...

Original headline:

blast off of apollo 8

Semantically closest:

blast off of apollo <u>13</u>

CSCS proposed change:

blast <u>victims</u> of apollo 8

# Some real example changes...

Original headline:

excuse me you left a gorilla suit on the bus


Semantically closest:

excuse me <u>we</u> left a gorilla suit on the bus


CSCS proposed change:

excuse me you left a gorilla <u>killer</u> on the bus

# Some real example changes...

Original headline:
winegrowers revel in good season

Semantically closest:
winegrowers revel in <u>strong</u> season

CSCS proposed change:
winegrowers revel in <u>flu</u> season

# Some real example changes...



Original headline:

winegrowers revel in good season

Semantically closest:

winegrowers revel in <u>strong</u> season

CSCS proposed change:

winegrowers revel in <u>flu</u> season

# A computational language model

On natural language sequences:

The American president _____
to Japan yesterday.

$p($"went"$) = 0.5$
$p($"traveled"$) = 0.2$
$p($"absconded"$) = 0.05$
$\vdots$
$p($"xylophone"$) = 0$

# A computational language model

On natural language sequences:

The American president _____
to Japan yesterday.

$p(\text{"went"}) = 0.5$
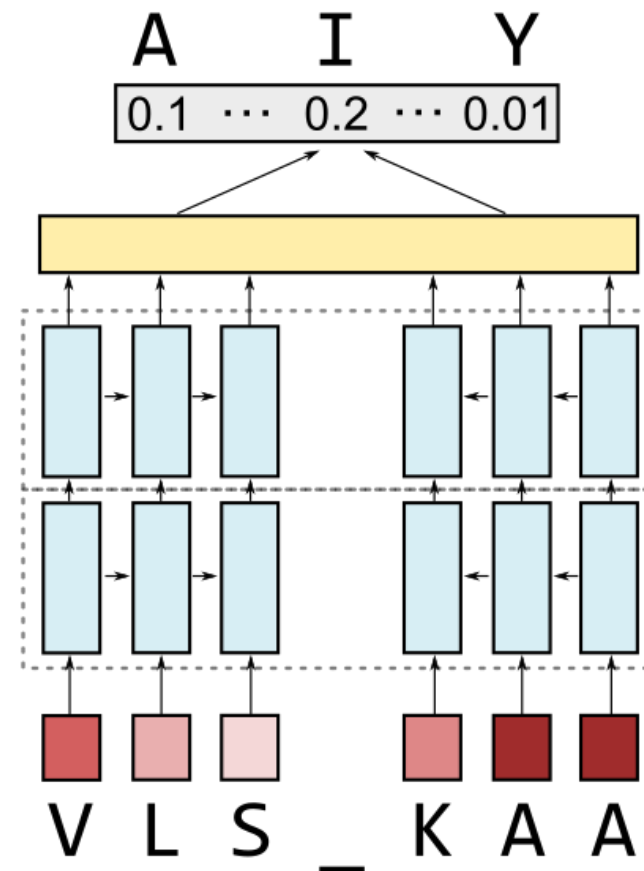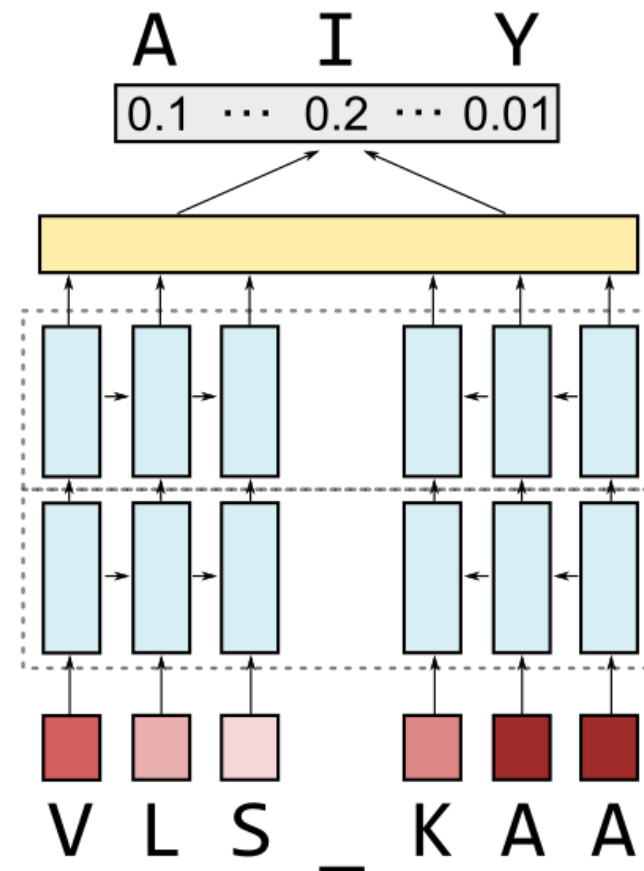$p(\text{"traveled"}) = 0.2$
$p(\text{"absconded"}) = 0.05$
　　　　$\vdots$
$p(\text{"xylophone"}) = 0$

# A computational language model

On natural language sequences:

The American president _____
to Japan yesterday.

$p(\text{"went"}) = 0.5$
$p(\text{"traveled"}) = 0.2$
$p(\text{"absconded"}) = 0.05$
     $\vdots$

$p(\text{"xyl}\ldots\text{"}) \ldots$

Trained on thousands of sequences (or more)!

# Train on viral protein sequence corpus



Influenza A HA
NIAID Influenza Research Database
(https://www.fludb.org)

HIV-1 Env
LANL HIV Database
(https://www.hiv.lanl.gov)

SARS-CoV-2 Spike
Virus Pathogen Database angd Analysis Resource
(https://www.viprbrc.org/)
GISAID
(https://www.gisaid.org/)
NCBI GenBank
(https://www.ncbi.nlm.nih.gov/sars-cov-2/)
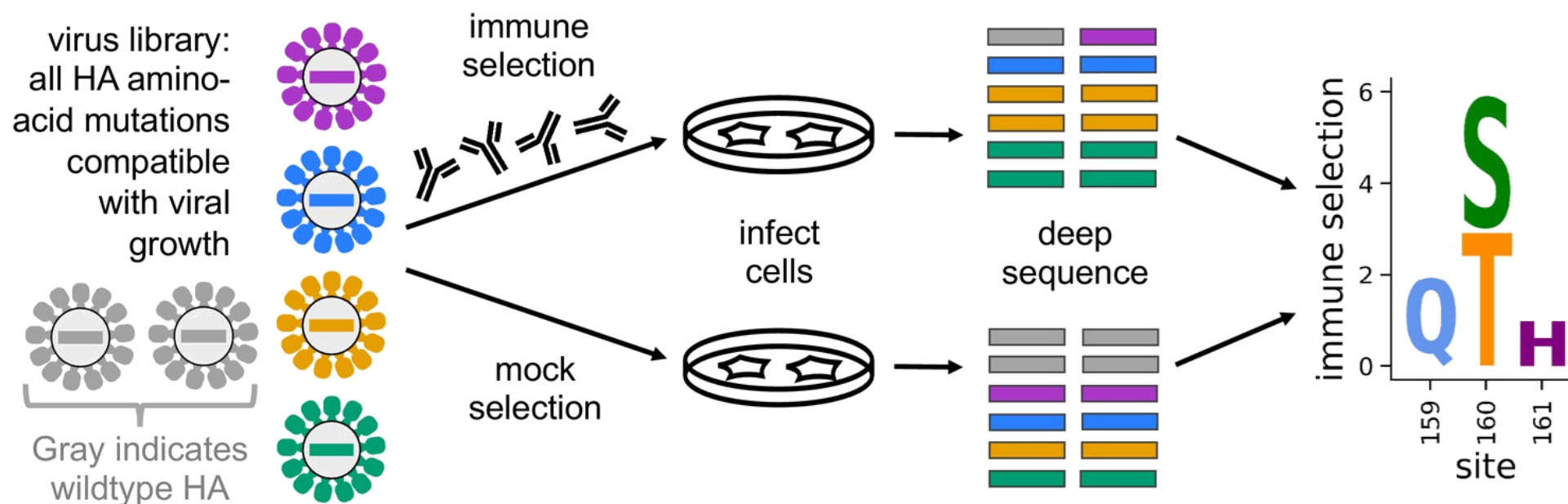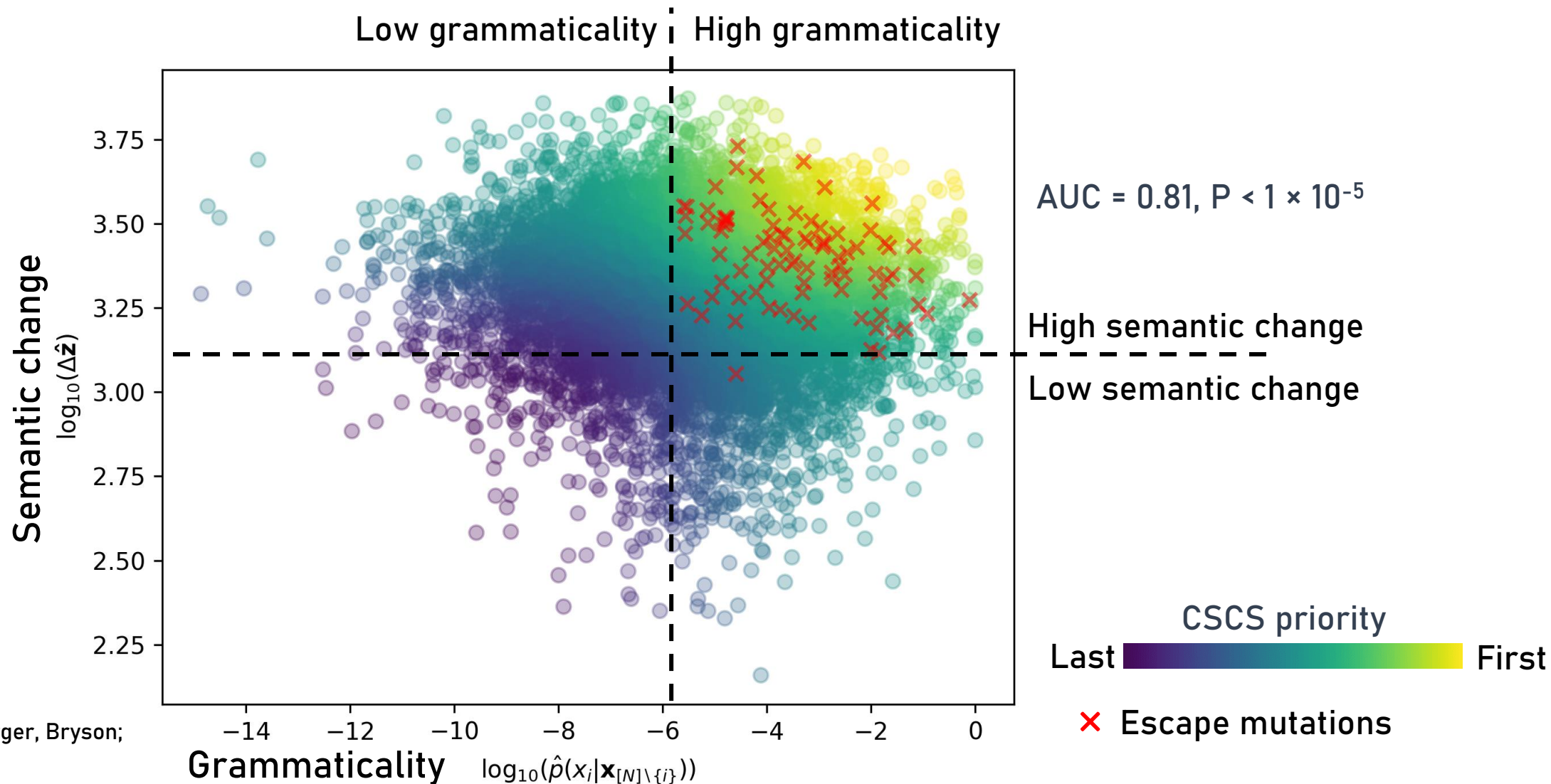
# Putting it all together to predict viral escape



From Lee et al., *eLife* (2019)

# Unsupervised prediction of escape mutations

Low grammaticality | High grammaticality

$\log_{10}(\Delta \hat{z})$ Semantic change

AUC = 0.81, P < 1 × 10$^{-5}$

High semantic change

Low semantic change

CSCS priority

Last ▬▬▬▬▬ First

× Escape mutations

Grammaticality    $\log_{10}(\hat{p}(x_i | \mathbf{x}_{[N] \setminus \{i\}}))$

From Hie, Zhong, Berger, Bryson;
*Science*, 2021

# Enriched escape potential in HA head



From Hie, Zhong, Berger, Bryson; *Science*, 2021

# Similar patterns for CoV-2 S1 versus S2



SARS-CoV-2 Spike

**N-terminal domain**
Escape enrichment,
$P < 1 \times 10^{-5}$

**Receptor binding domain**
Escape enrichment,
$P = 2.7 \times 10^{-3}$

**RBD**

90°

**S2**
Escape depletion,
$P < 1 \times 10^{-5}$

Predicted
escape potential
– +

From Hie, Zhong, Berger, Bryson; *Science*, 2021

# Language model predicts SARS-CoV-2 variants



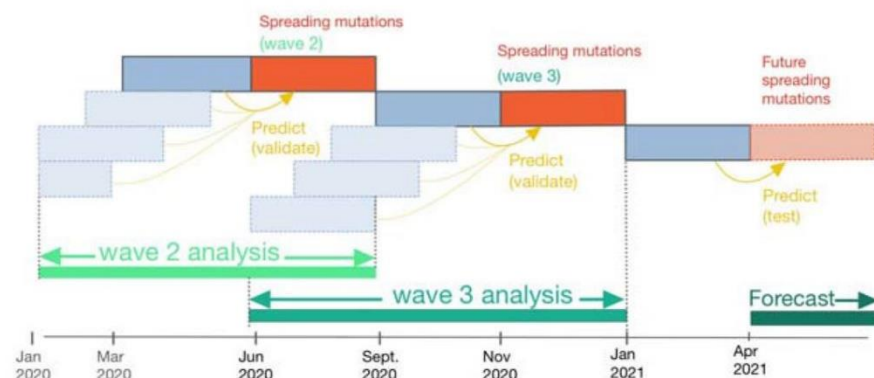Spike language model predicts mutations up to 4 months in advance with AUC of 0.8

Cyrus Maher

# Language model predicts SARS-CoV-2 variants
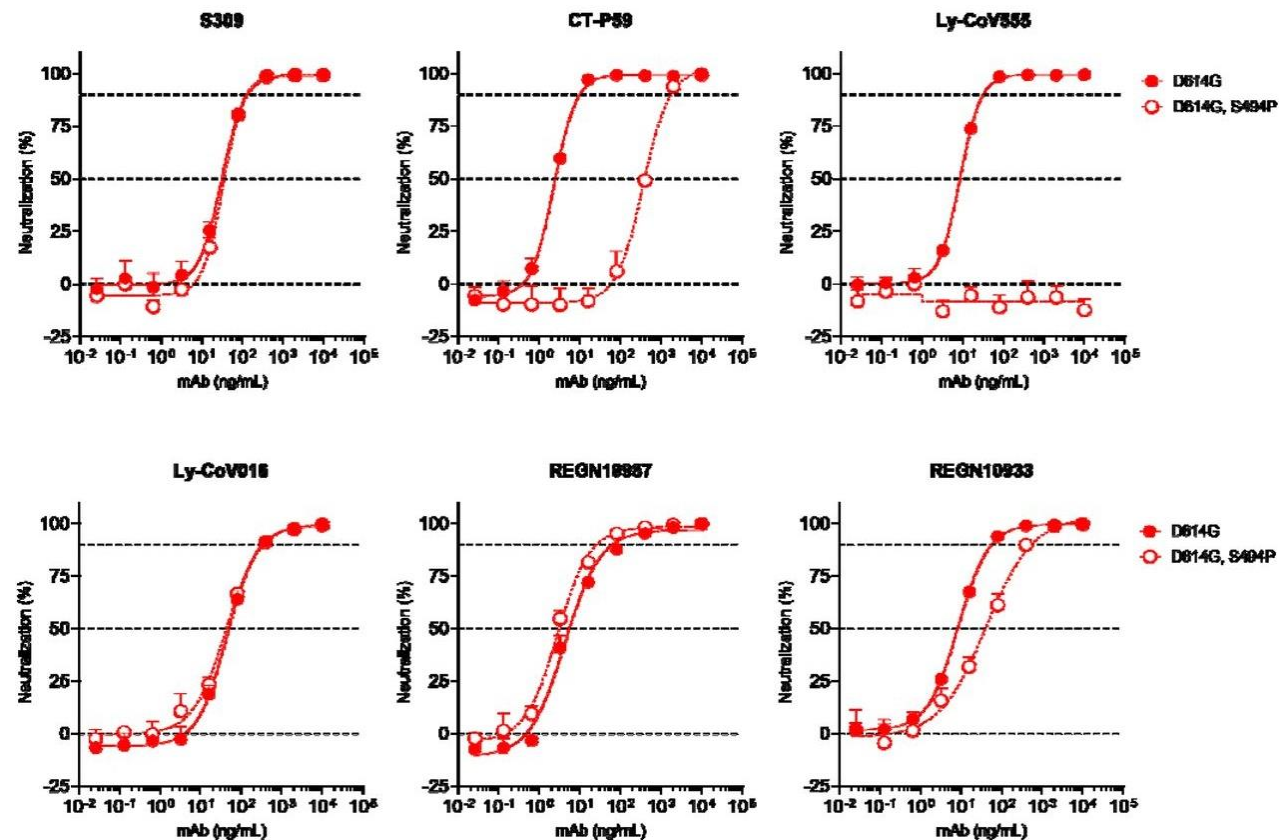


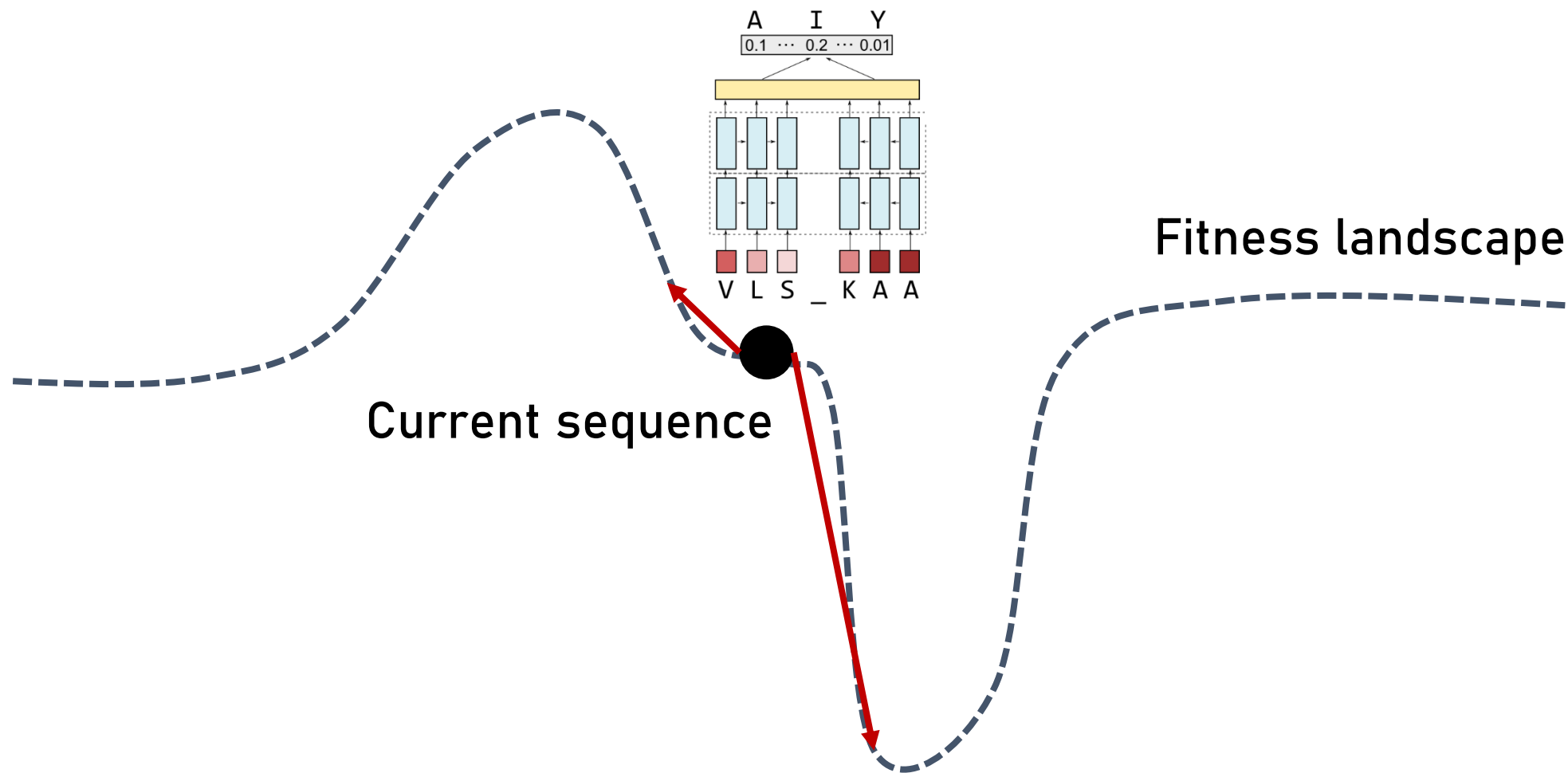Spike language model predicts mutations up to 4 months in advance with AUC of 0.8

Cyrus Maher

From Maher et al.; *medRxiv*, 2021

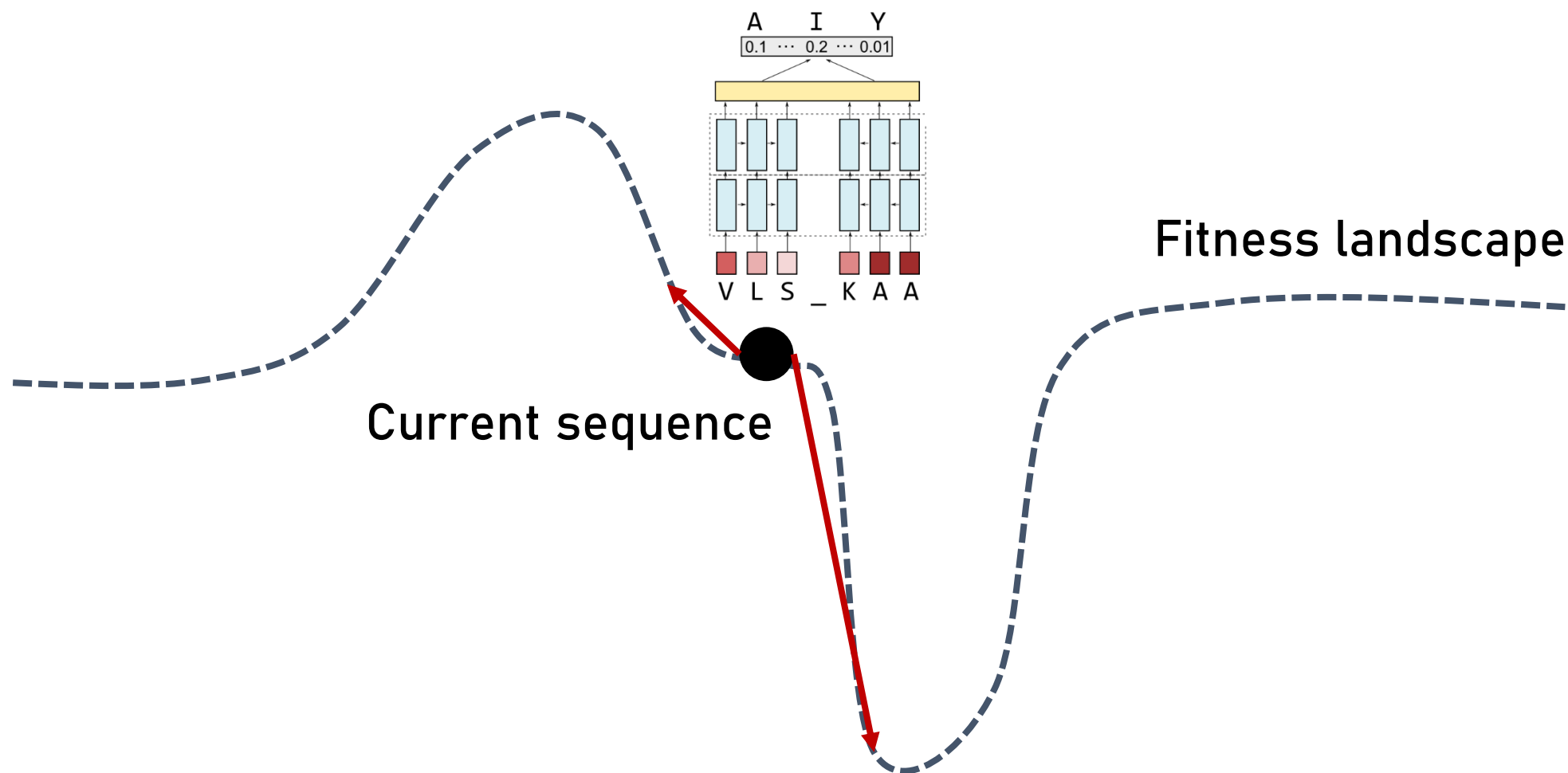1. Learning the language of viral evolution and escape

2. Recovering evolutionary dynamics with "evolutionary velocity"

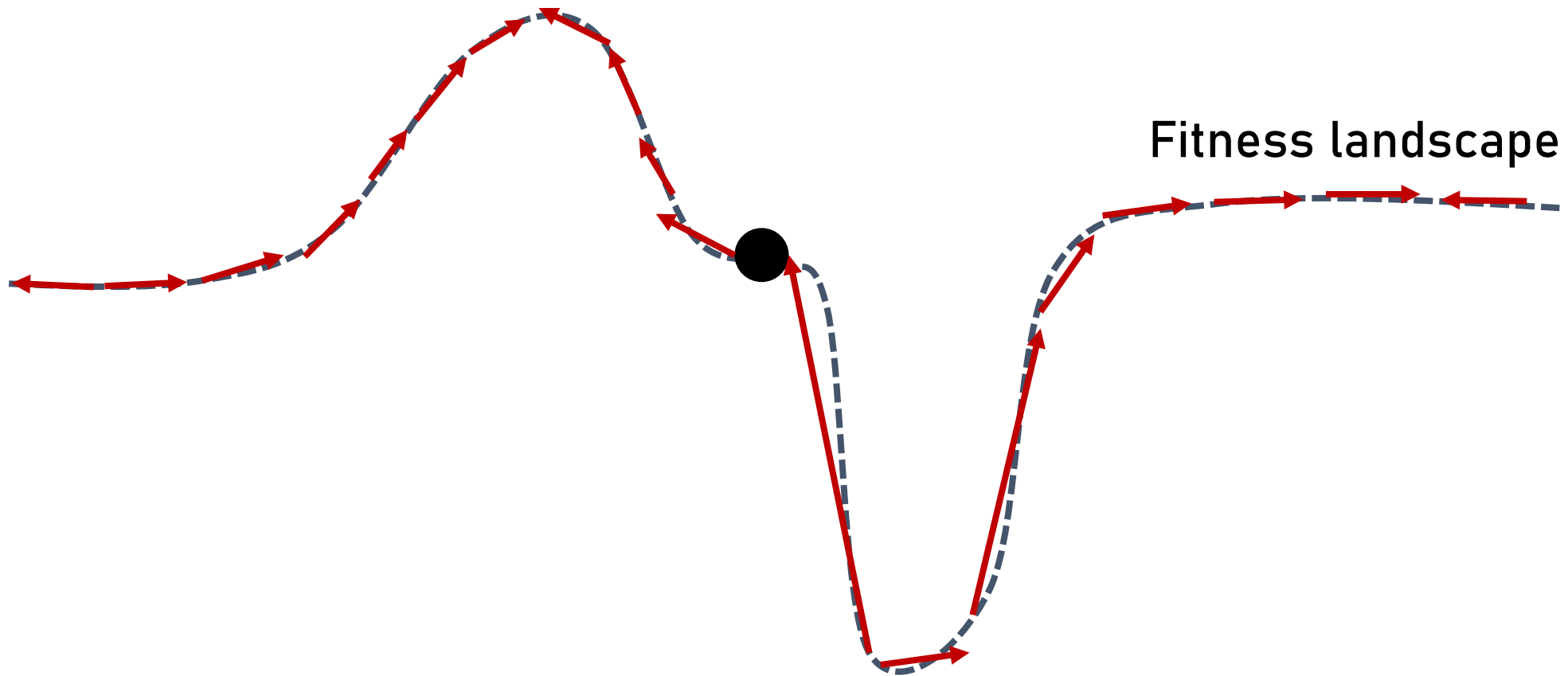3. Looking forward

# We can predict local evolution with LMs



Fitness landscape

Current sequence

# Understand <u>global</u> patterns using <u>local</u> predictions
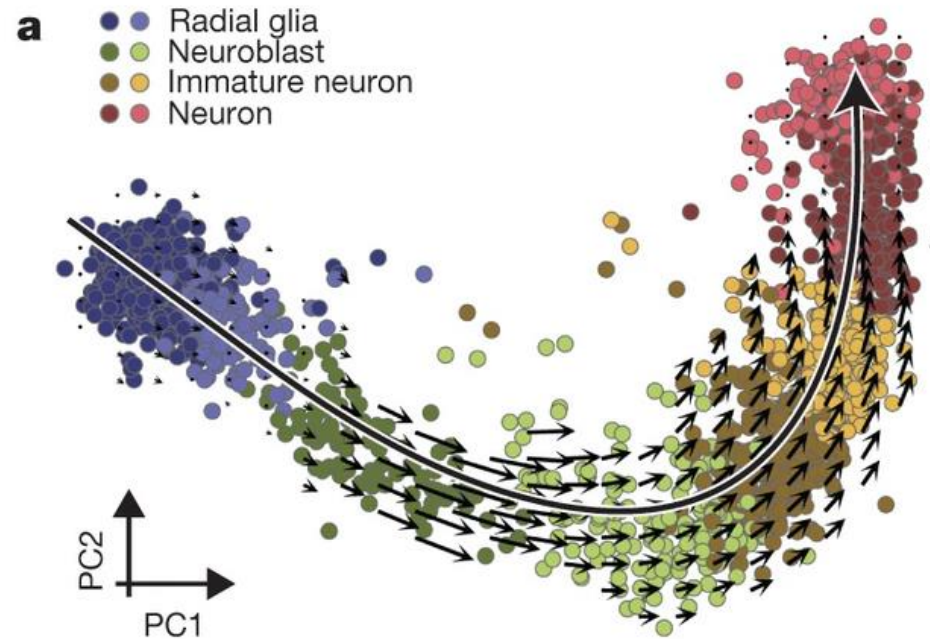


Fitness landscape

# Understand global patterns using local predictions

"RNA velocity"
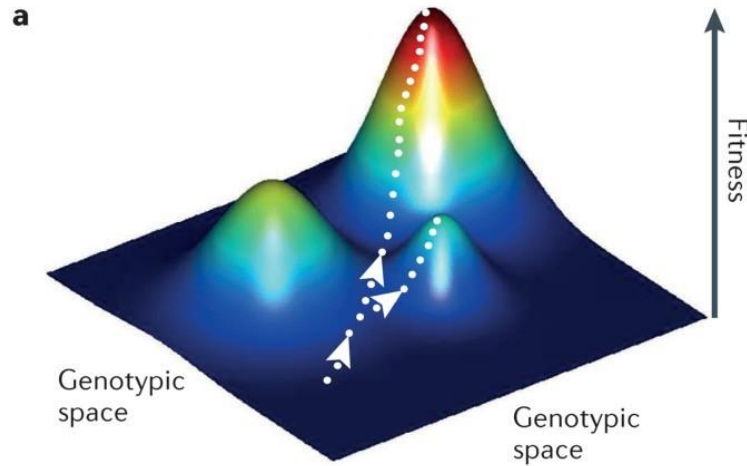
(La Manno et al, *Nature*, 2018)

# Understand <u>global</u> patterns using <u>local</u> predictions
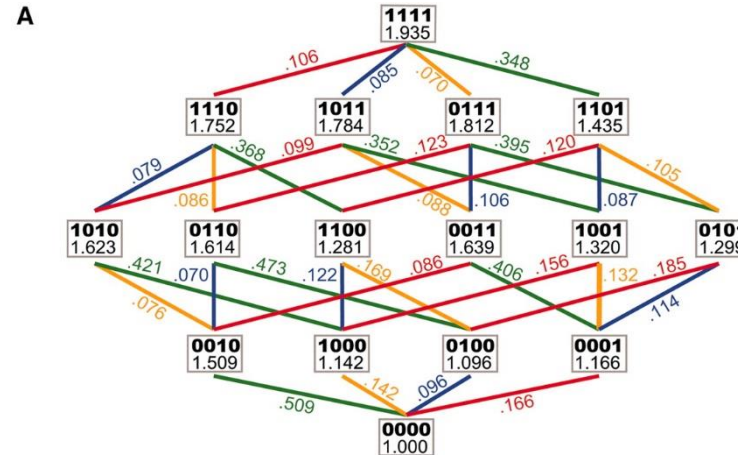


"Fitness landscape"
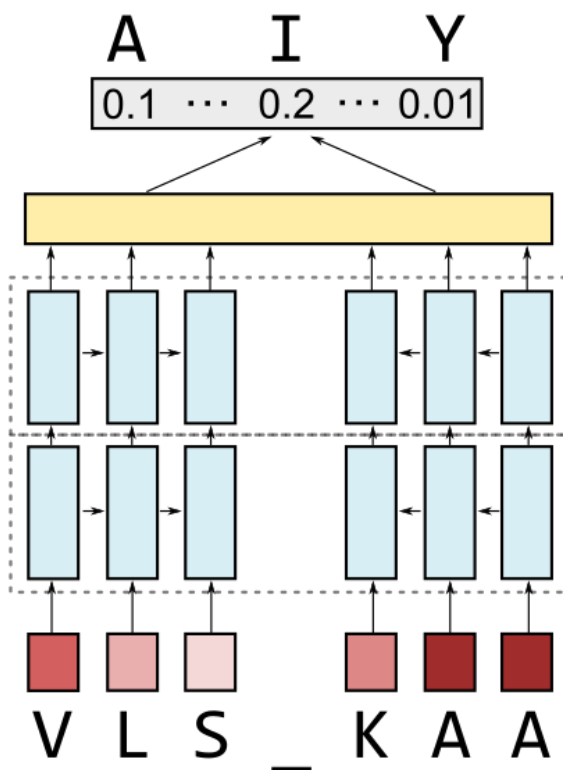
(Wright, *Int. Conf. Genetics*, 1932)



Visser and Krug. *Nat. Rev. Genetics* (2014)
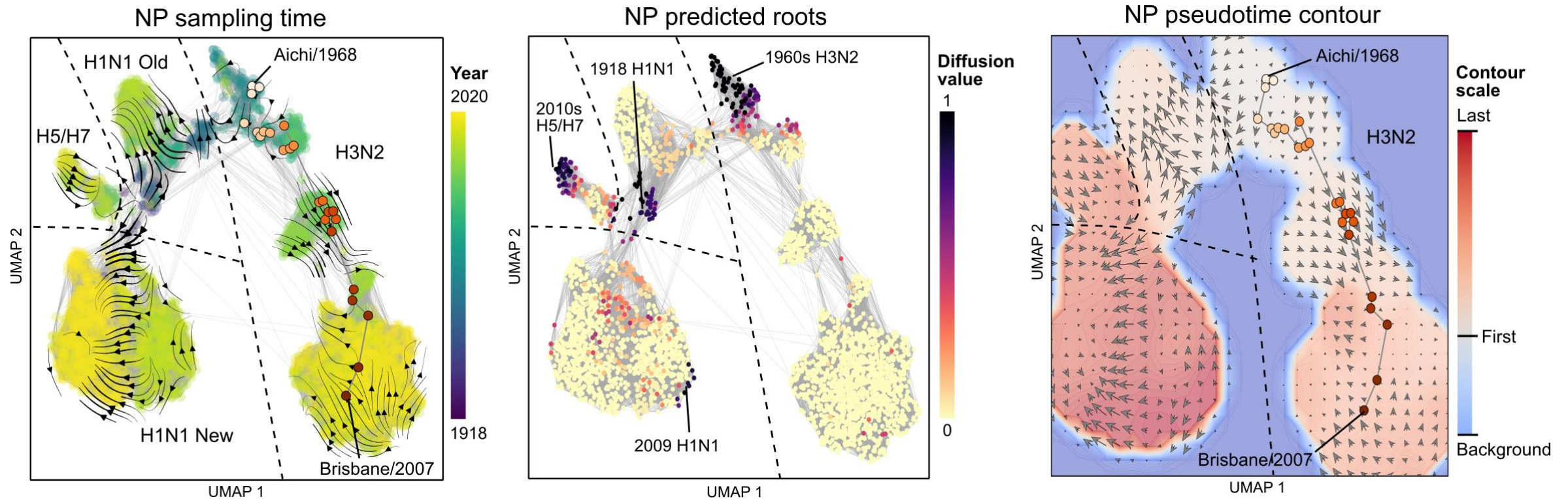
Chou et al. *Science* (2011)

# "Universal" protein language model
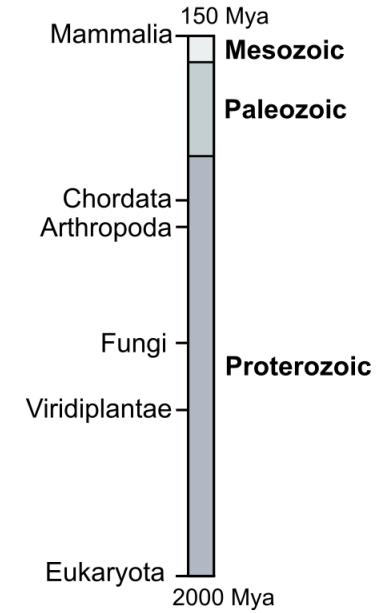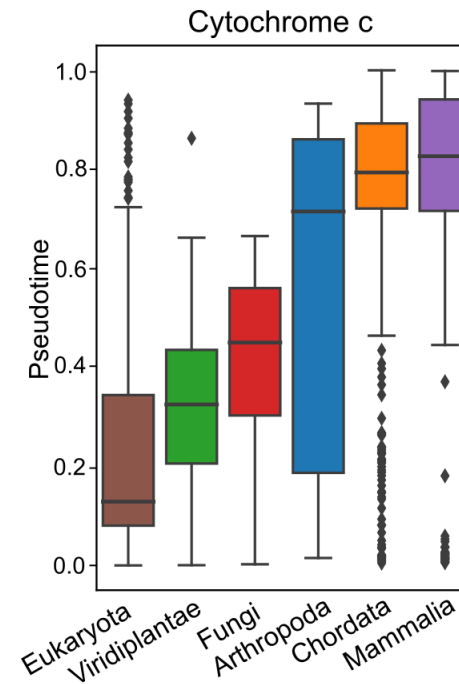

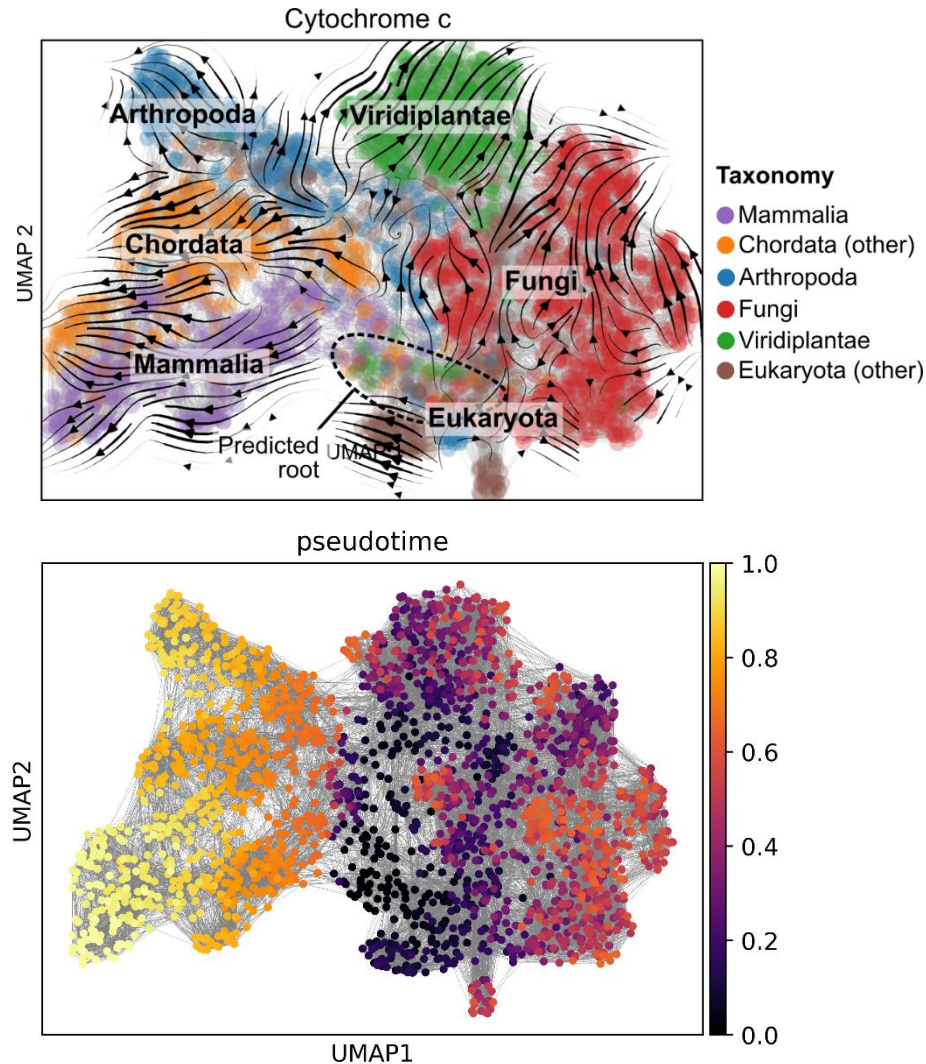
ESM-1b by Rives et al. (*PNAS*, 2021):
- Trained on 3 million sequences from UniRef50
- Model has 650 million parameters

# Velocity of influenza evolution



NP sampling time

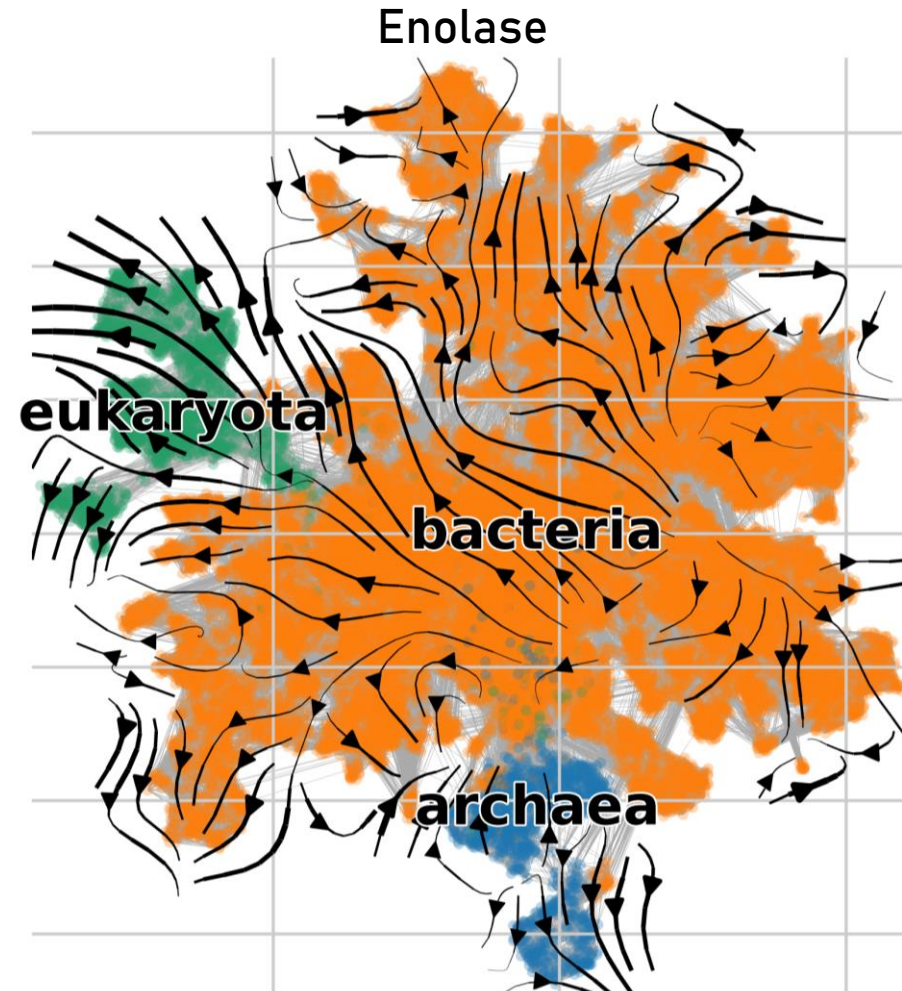NP predicted roots

NP pseudotime contour

Temporal Spearman r = 0.49, P < 1e-308

From Hie, Yang, and Kim; *bioRxiv*, 2021

# Velocity of cytochrome c evolution



From Hie, Yang, and Kim; *bioRxiv*, 2021

# Velocity of ancient evolution



PGK

Enolase
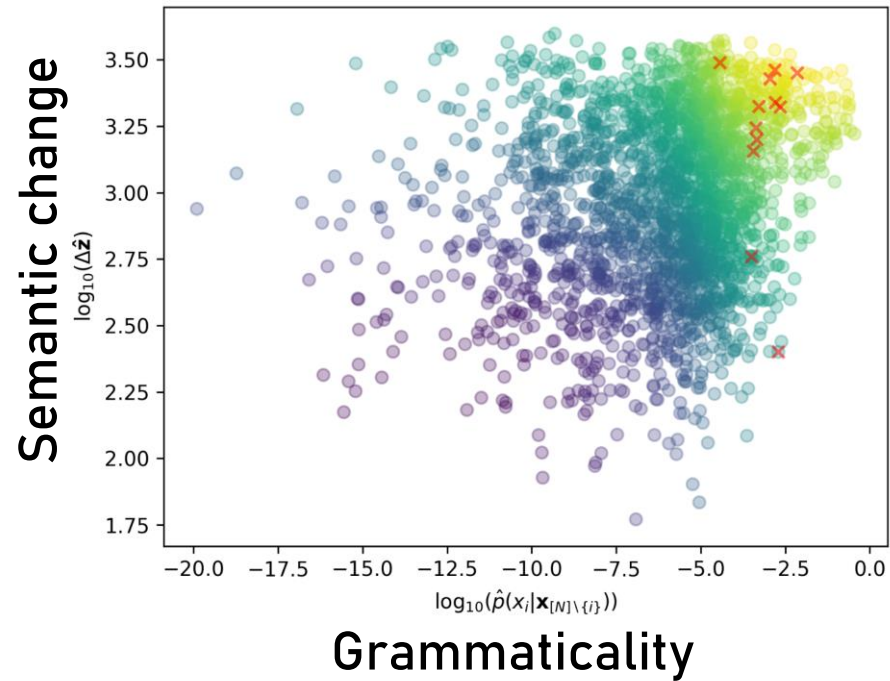
From Hie, Yang, and Kim; *bioRxiv*, 2021

1. Learning the language of viral evolution and escape

2. Recovering evolutionary dynamics with "evolutionary velocity"

3. Looking forward

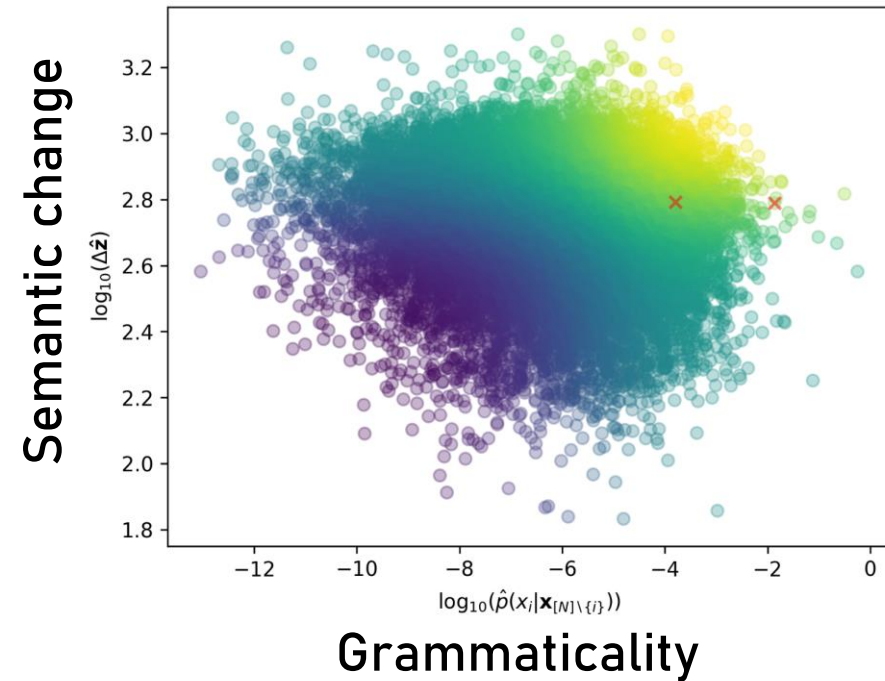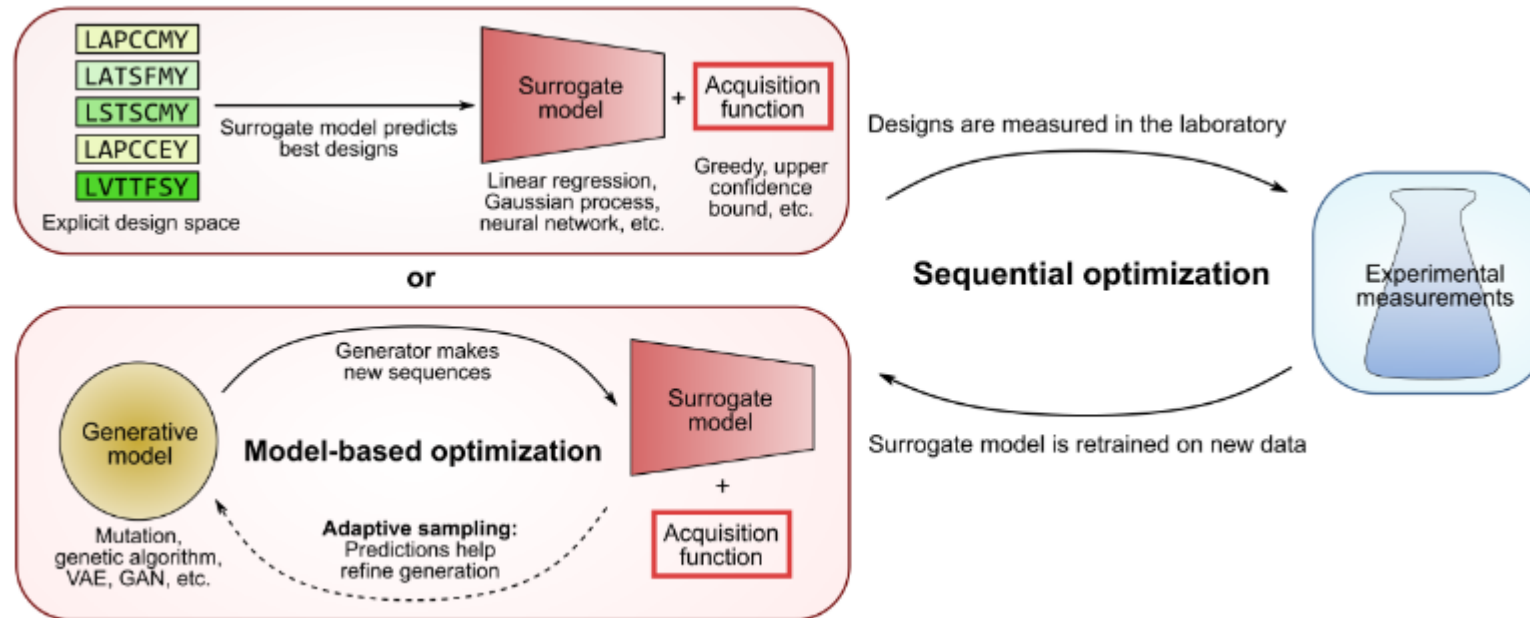# Predicting drug resistance



HIV

M. tuberculosis

Andy Tso

# Evolving protein therapeutics



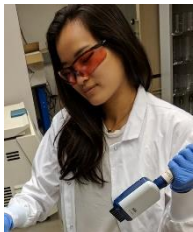From Hie and Yang; *arXiv*, 2021

# Key takeaways

- Language models have the potential to improve evolutionary models and prediction

- Sufficient training data is important!

- Successful implementation will require interdisciplinary collaboration

# References

- Hie, Zhong, Berger, and Bryson.
  "Learning the language of viral evolution and escape."
  *Science,* 371:6526 (2021).

- Maher, ..., Hie et al.
  "Predicting the mutational drivers of future SARS-CoV-2 variants of concern."
  *medRxiv* (2021).

- Hie, Yang, and Kim.
  "Evolutionary velocity with protein language models."
  *bioRxiv* (2021).

# Thank you!



Ellen Zhong

Bryan Bryson

Bonnie Berger

Kevin Yang

Peter Kim