



# Protein Language Models

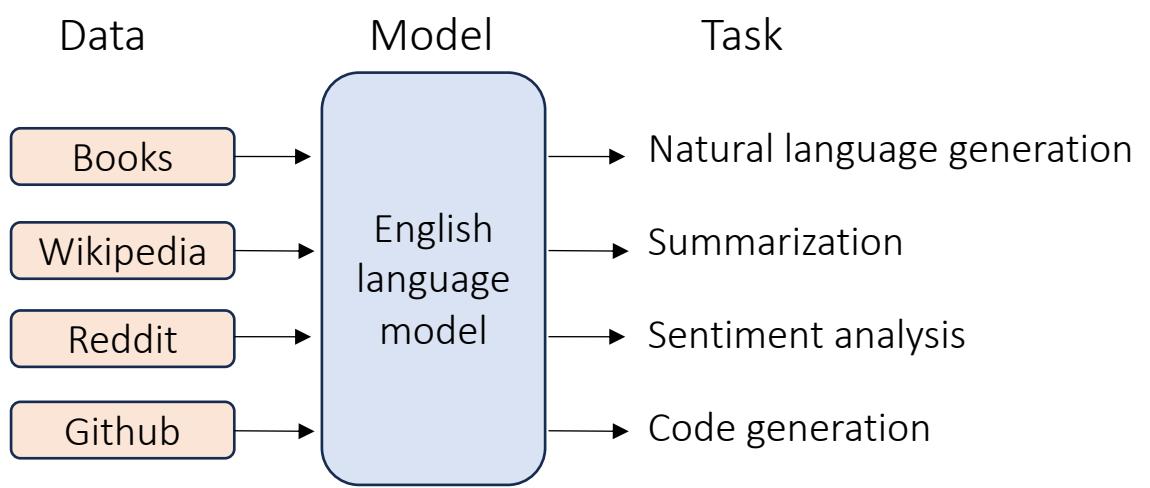
Sarah Gurev  
MLCB 2024

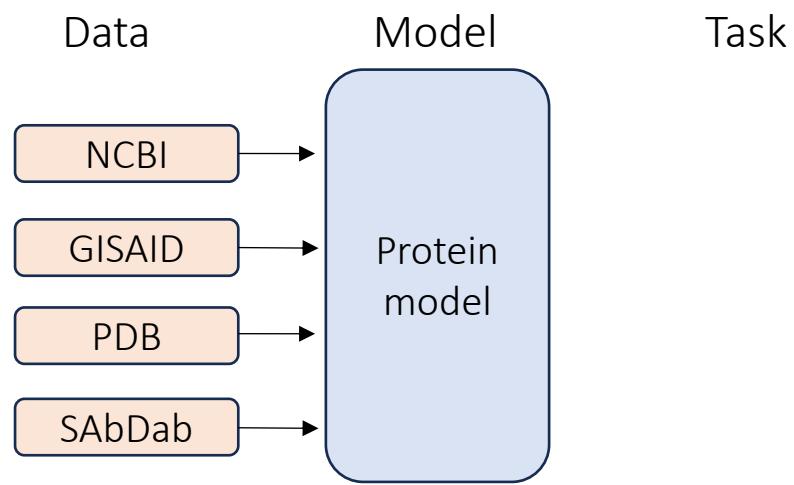
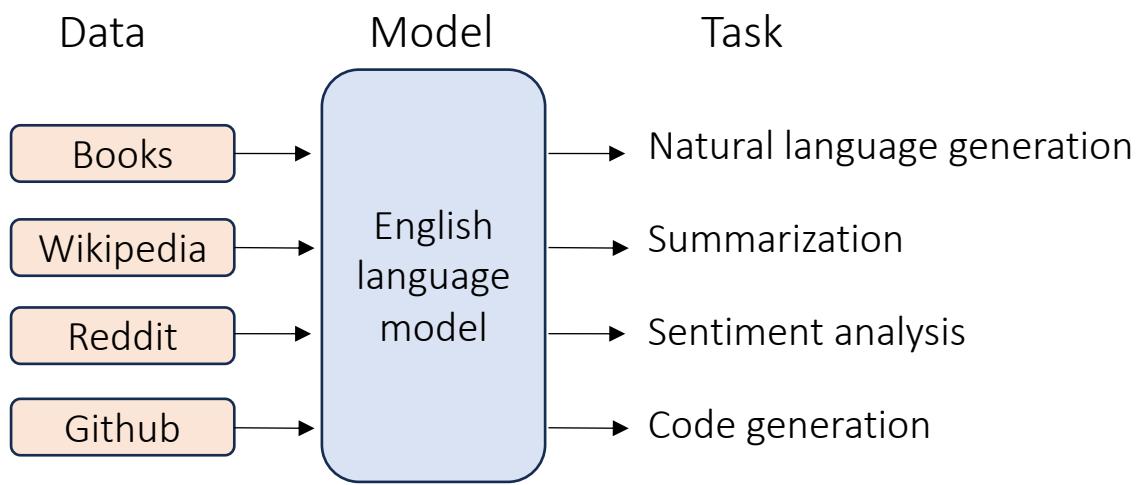
# Parallels between natural languages & proteins

| DNA / RNA                         | Protein               |
|-----------------------------------|-----------------------|
| <i>4 nucleotides or 64 codons</i> | <i>20 amino acids</i> |
| ATGTTCATCGTCCTG...                | MADRLYMTKIHHQFDGD...  |

# Parallels between natural languages & proteins

| Languages          | DNA / RNA                         | Protein               |
|--------------------|-----------------------------------|-----------------------|
| <i>26 alphabet</i> | <i>4 nucleotides or 64 codons</i> | <i>20 amino acids</i> |
| <b>String</b>      | The cat sat on the mat.           | ATGTTCATCGTCCTG...    |

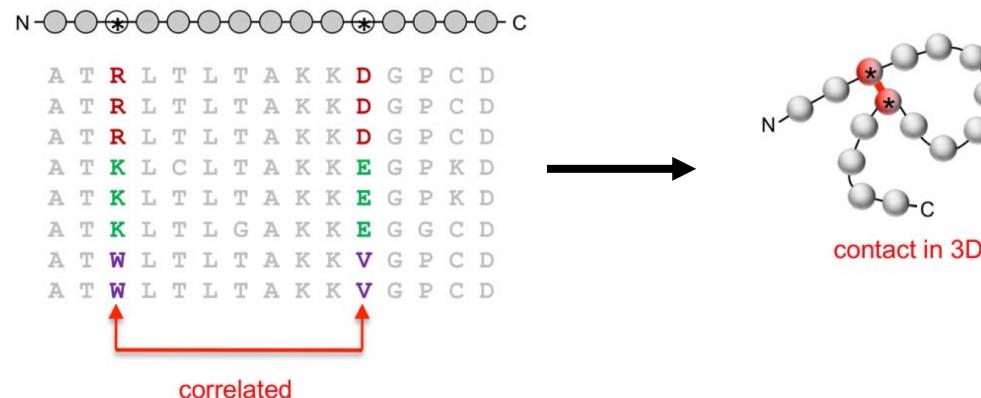




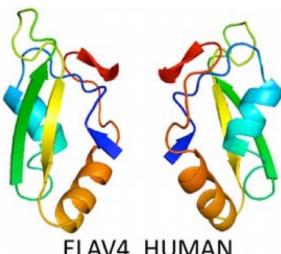
# Protein 3D Structure Computed from Evolutionary Sequence Variation

Debora S. Marks<sup>1\*3</sup>, Lucy J. Colwell<sup>2\*</sup>, Robert Sheridan<sup>3</sup>, Thomas A. Hopf<sup>1</sup>, Andrea Pagnani<sup>4</sup>, Riccardo Zecchina<sup>4,5</sup>, Chris Sander<sup>3</sup>

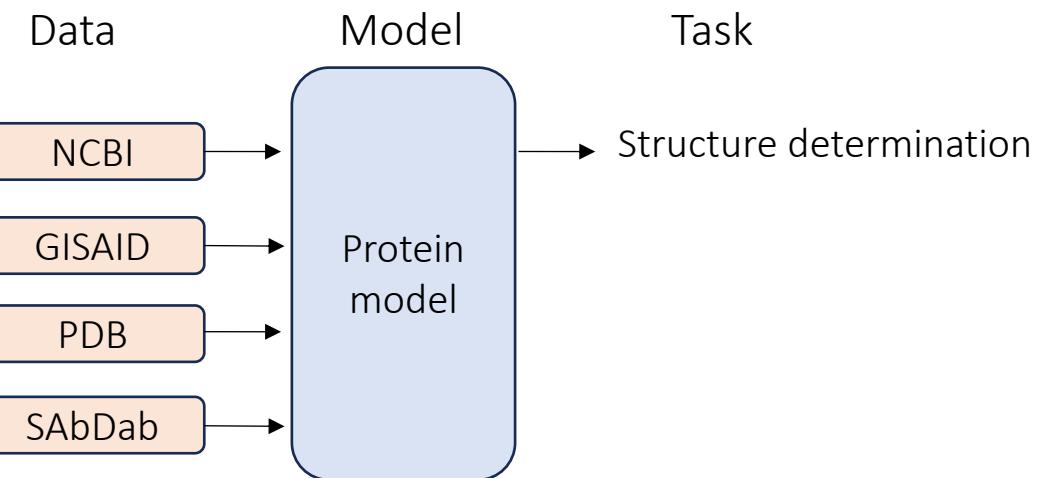
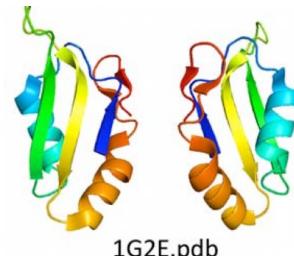
**1** Department of Systems Biology, Harvard Medical School, Boston, Massachusetts, United States of America, **2** MRC Laboratory of Molecular Biology, Hills Road, Cambridge, United Kingdom, **3** Computational Biology Center, Memorial Sloan-Kettering Cancer Center, New York, New York, United States of America, **4** Human Genetics Foundation, Torino, Italy, **5** Politecnico di Torino, Torino, Italy



Predicted



Observed



# Sequence co-evolution gives 3D contacts and structures of protein complexes

Thomas A Hopf<sup>1,2†</sup>, Charlotta P I Schärfe<sup>1,3,4†</sup>, João P G L M Rodrigues<sup>5†</sup>,  
Anna G Green<sup>1</sup>, Oliver Kohlbacher<sup>3,4</sup>, Chris Sander<sup>6\*</sup>, Alexandre M J J Bonvin<sup>5\*</sup>,  
Debora S Marks<sup>1\*</sup>

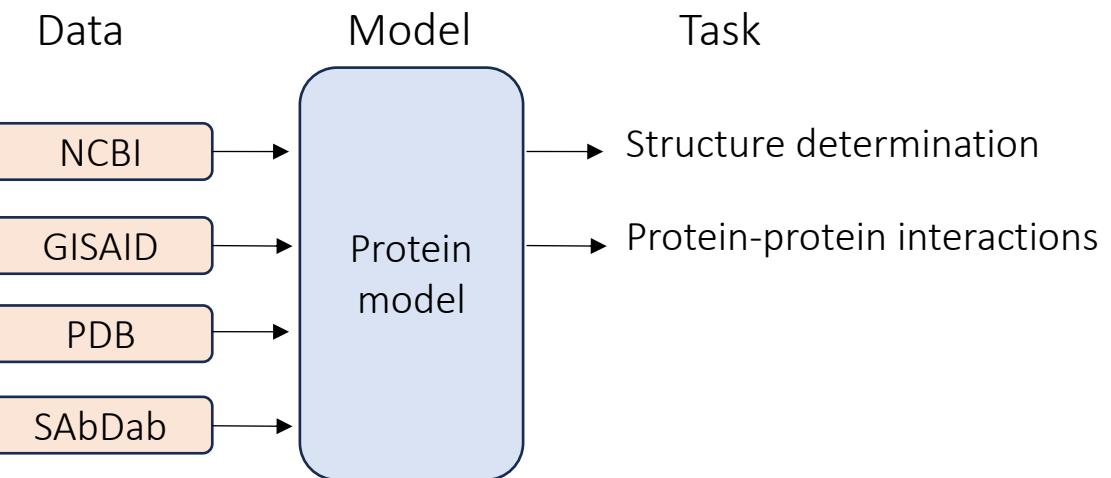
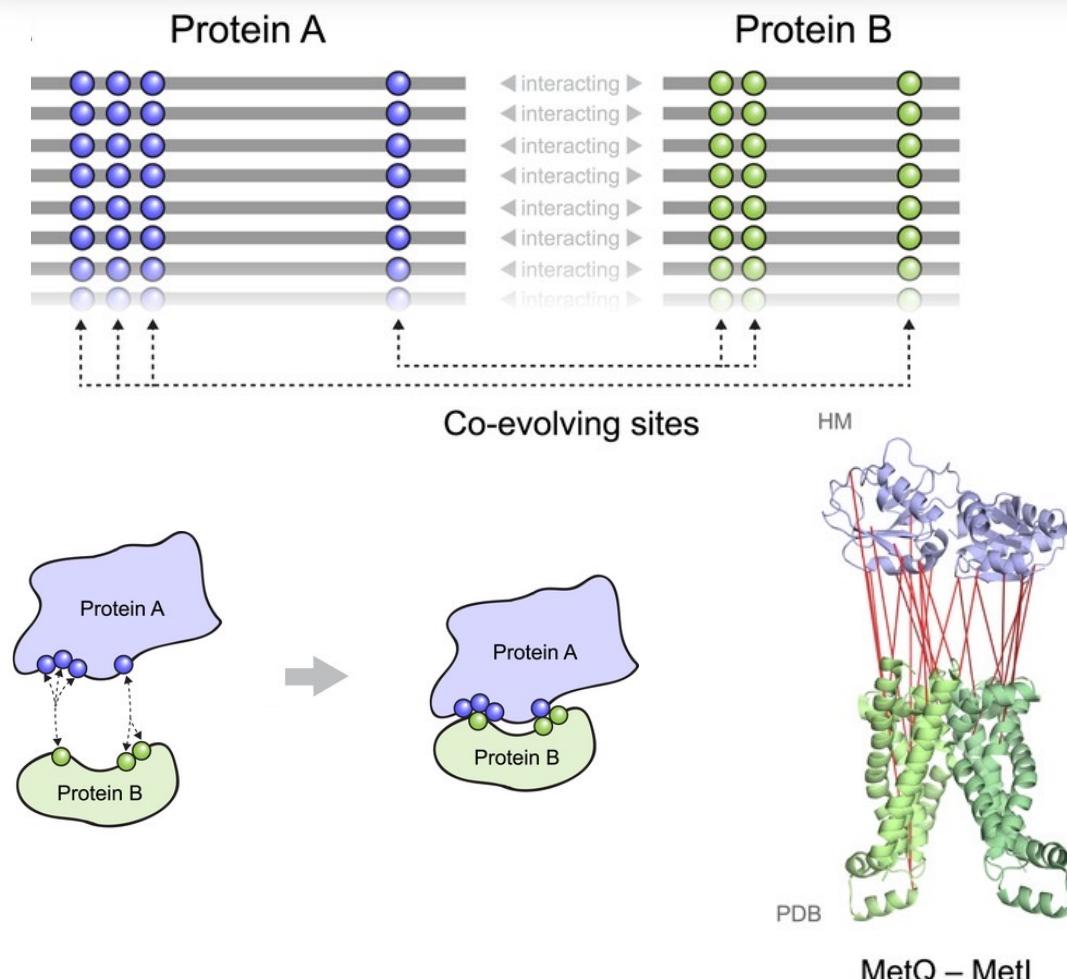


Figure adapted from Madani et al., 2023. Nature Biotech

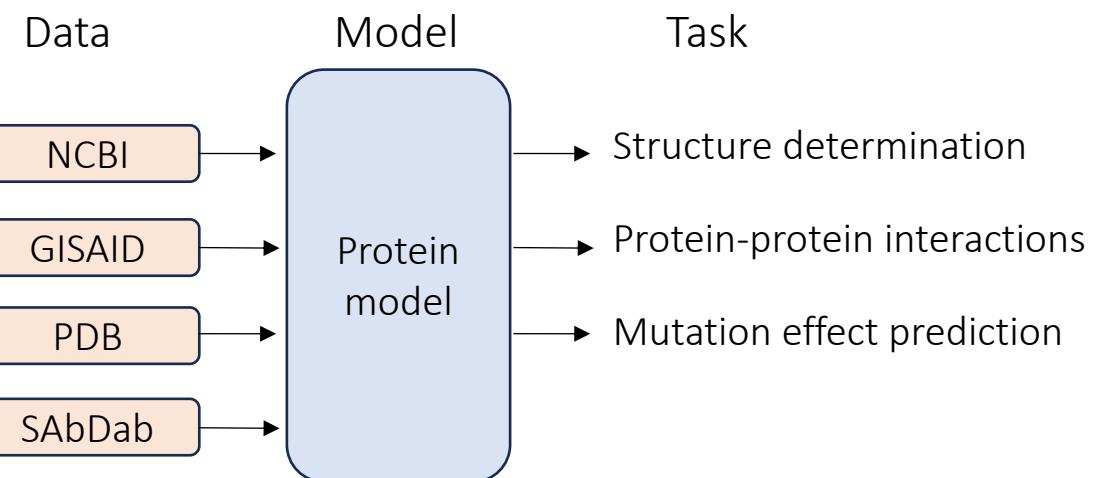
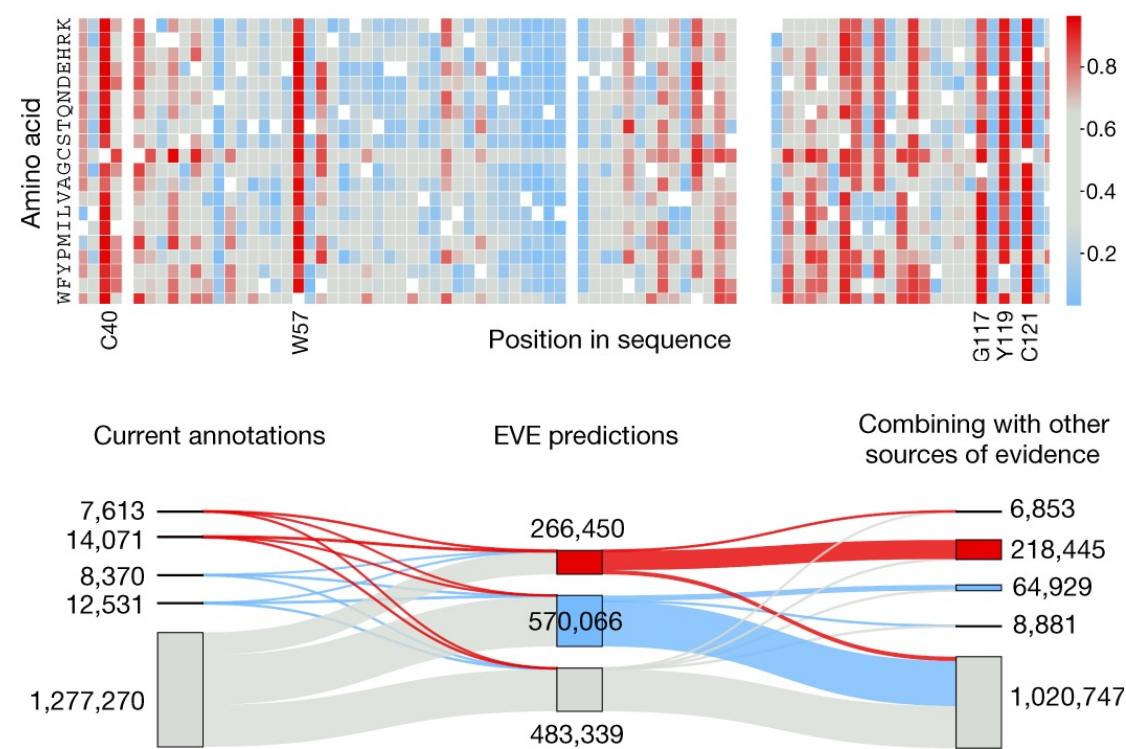
## Article

# Disease variant prediction with deep generative models of evolutionary data

<https://doi.org/10.1038/s41586-021-04043-8>

Jonathan Frazer<sup>1,4</sup>, Pascal Notin<sup>2,4</sup>, Mafalda Dias<sup>1,4</sup>, Aidan Gomez<sup>2</sup>, Joseph K. Min<sup>1</sup>, Kelly Brock<sup>1</sup>, Yarin Gal<sup>2,3,✉</sup> & Debora S. Marks<sup>1,3,✉</sup>

Received: 18 December 2020



Article

# De novo design of protein structure and function with RFdiffusion

<https://doi.org/10.1038/s41586-023-06415-8>

Received: 14 December 2022

Accepted: 7 July 2023

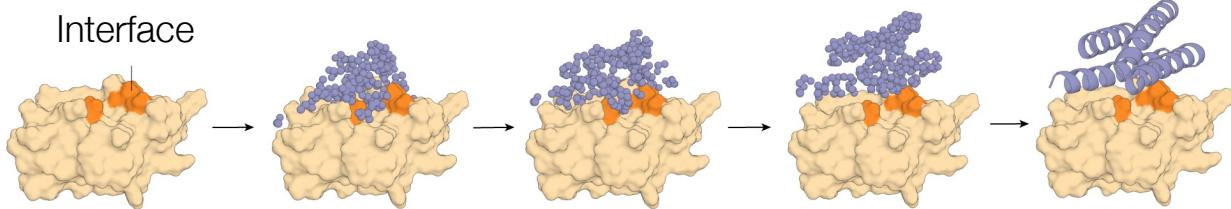
Published online: 11 July 2023

Open access

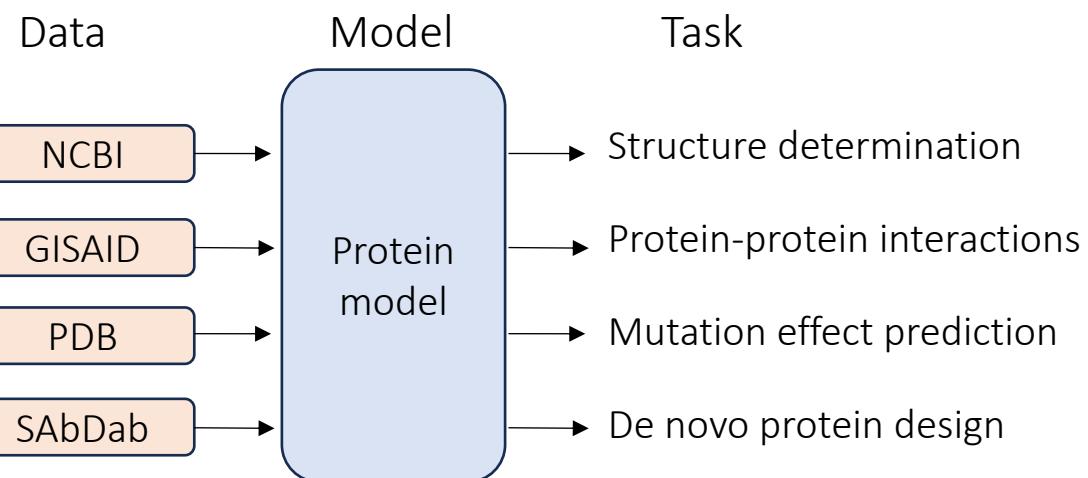
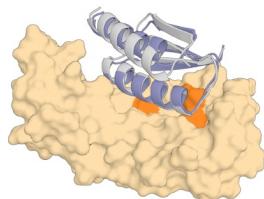
 Check for updates

Joseph L. Watson<sup>1,2,15</sup>, David Juergens<sup>1,2,3,15</sup>, Nathaniel R. Bennett<sup>1,2,3,15</sup>, Brian L. Tripp<sup>2,4,5,15</sup>, Jason Yim<sup>2,6,15</sup>, Helen E. Eisenach<sup>1,2,15</sup>, Woody Ahern<sup>1,2,7,15</sup>, Andrew J. Borst<sup>1,2</sup>, Robert J. Ragotte<sup>1,2</sup>, Lukas F. Milles<sup>1,2</sup>, Basile I. M. Wicky<sup>1,2</sup>, Nikita Hanikel<sup>1,2</sup>, Samuel J. Pellock<sup>1,2</sup>, Alexis Courbet<sup>1,2,8</sup>, William Sheffler<sup>1,2</sup>, Jue Wang<sup>1,2</sup>, Preetham Venkatesh<sup>1,2,9</sup>, Isaac Sappington<sup>1,2,9</sup>, Susana Vázquez Torres<sup>1,2,9</sup>, Anna Lauko<sup>1,2,9</sup>, Valentin De Bortoli<sup>8</sup>, Emile Mathieu<sup>10</sup>, Sergey Ovchinnikov<sup>11,12</sup>, Regina Barzilay<sup>6</sup>, Tommi S. Jaakkola<sup>6</sup>, Frank DiMaio<sup>1,2</sup>, Minkyung Baek<sup>13</sup> & David Baker<sup>1,2,14</sup>✉

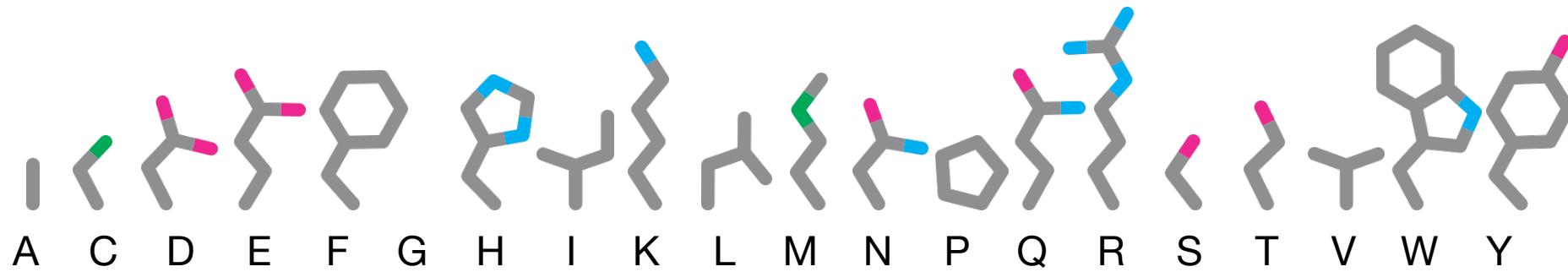
## Interface



Novel binder to influenza HA protein

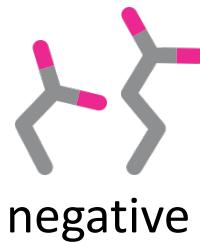


# The protein alphabet

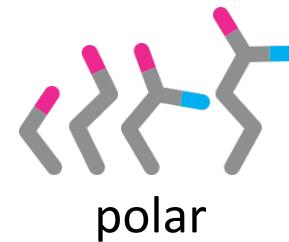




positive



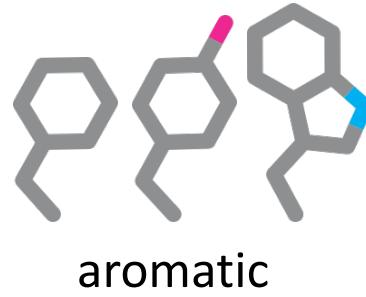
negative



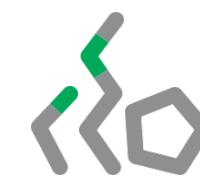
polar



"greasy"



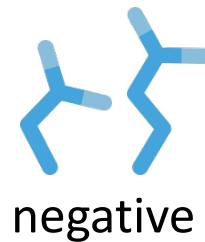
aromatic



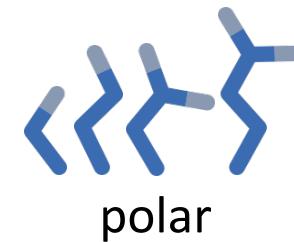
other stuff



positive



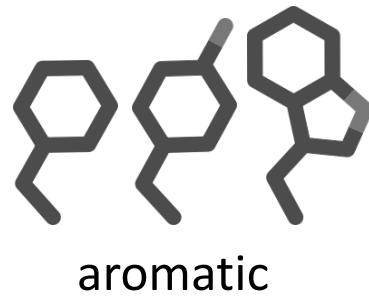
negative



polar



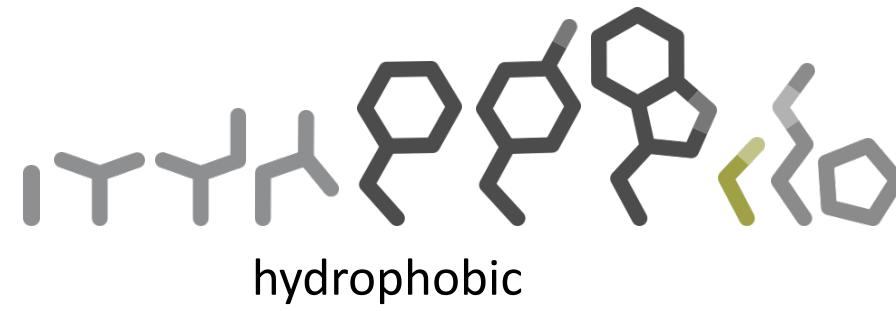
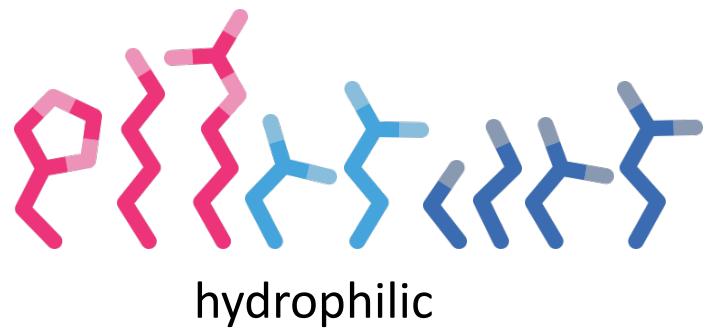
“greasy”

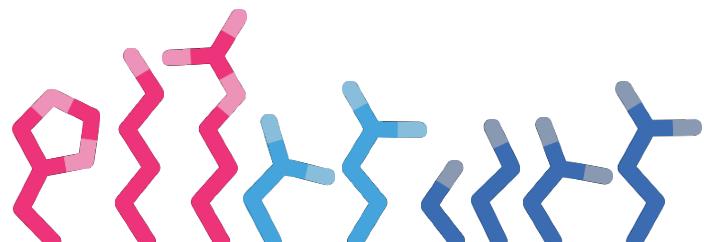


aromatic

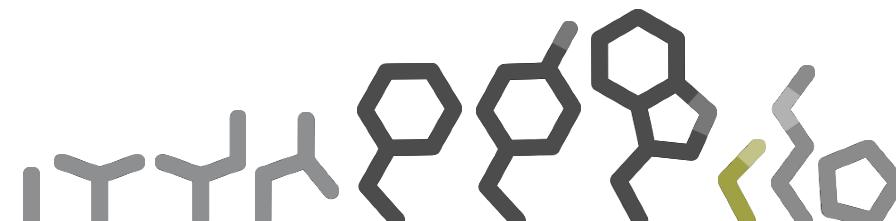


other stuff

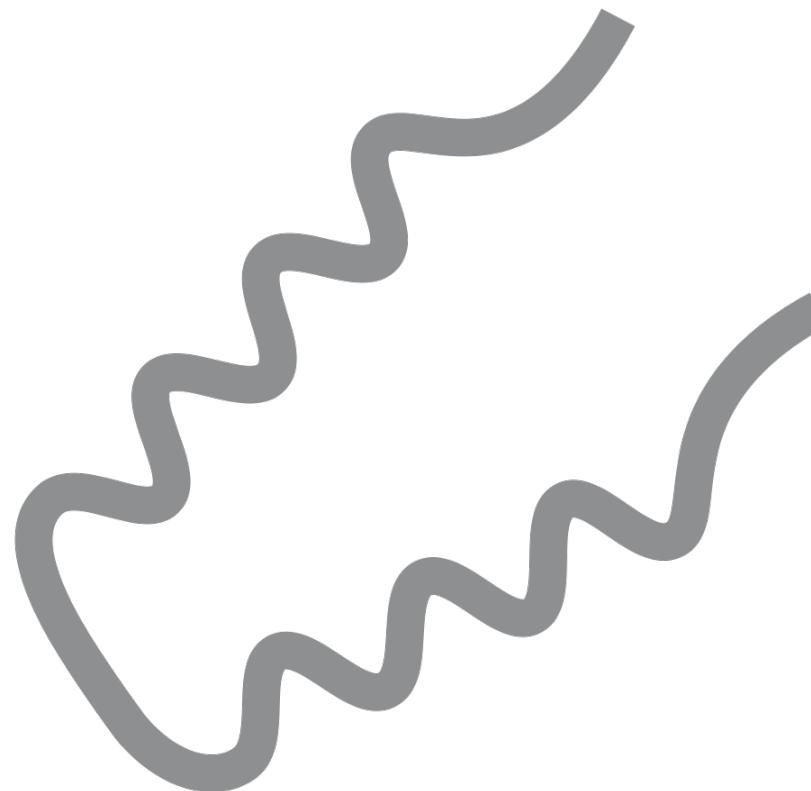


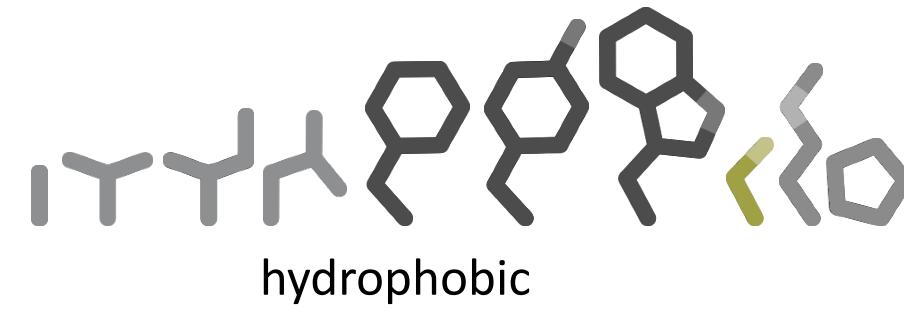
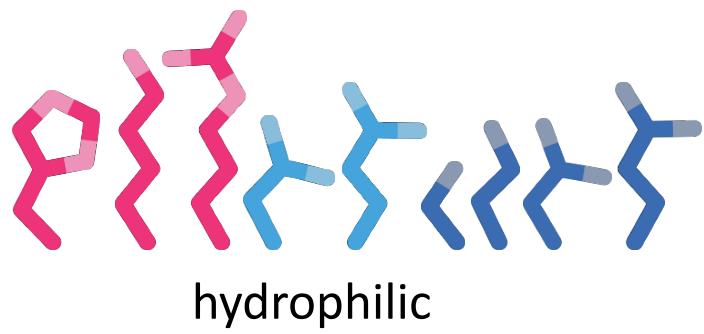
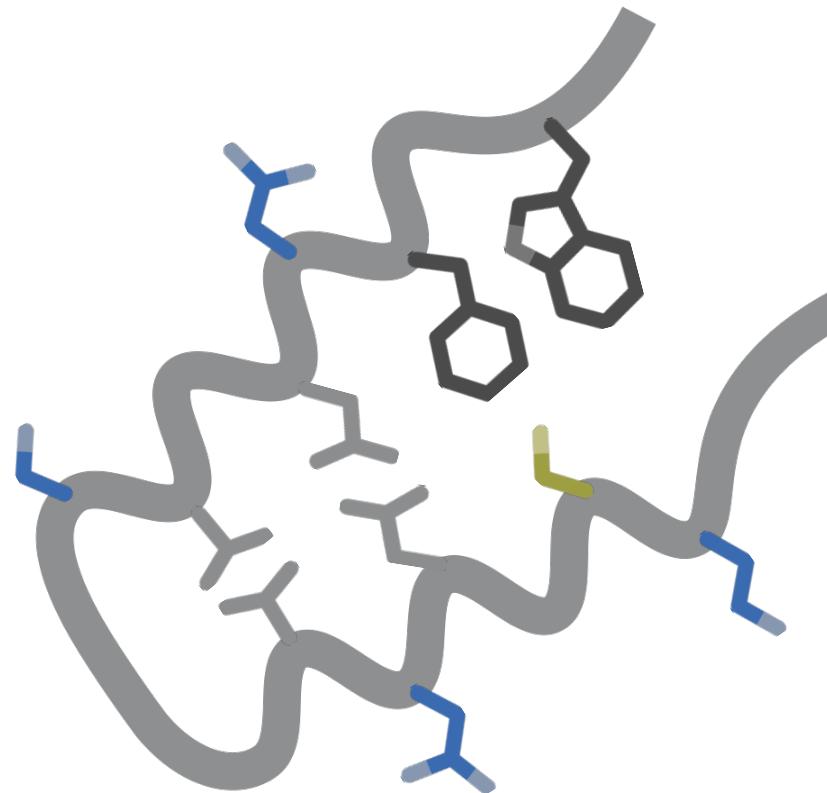


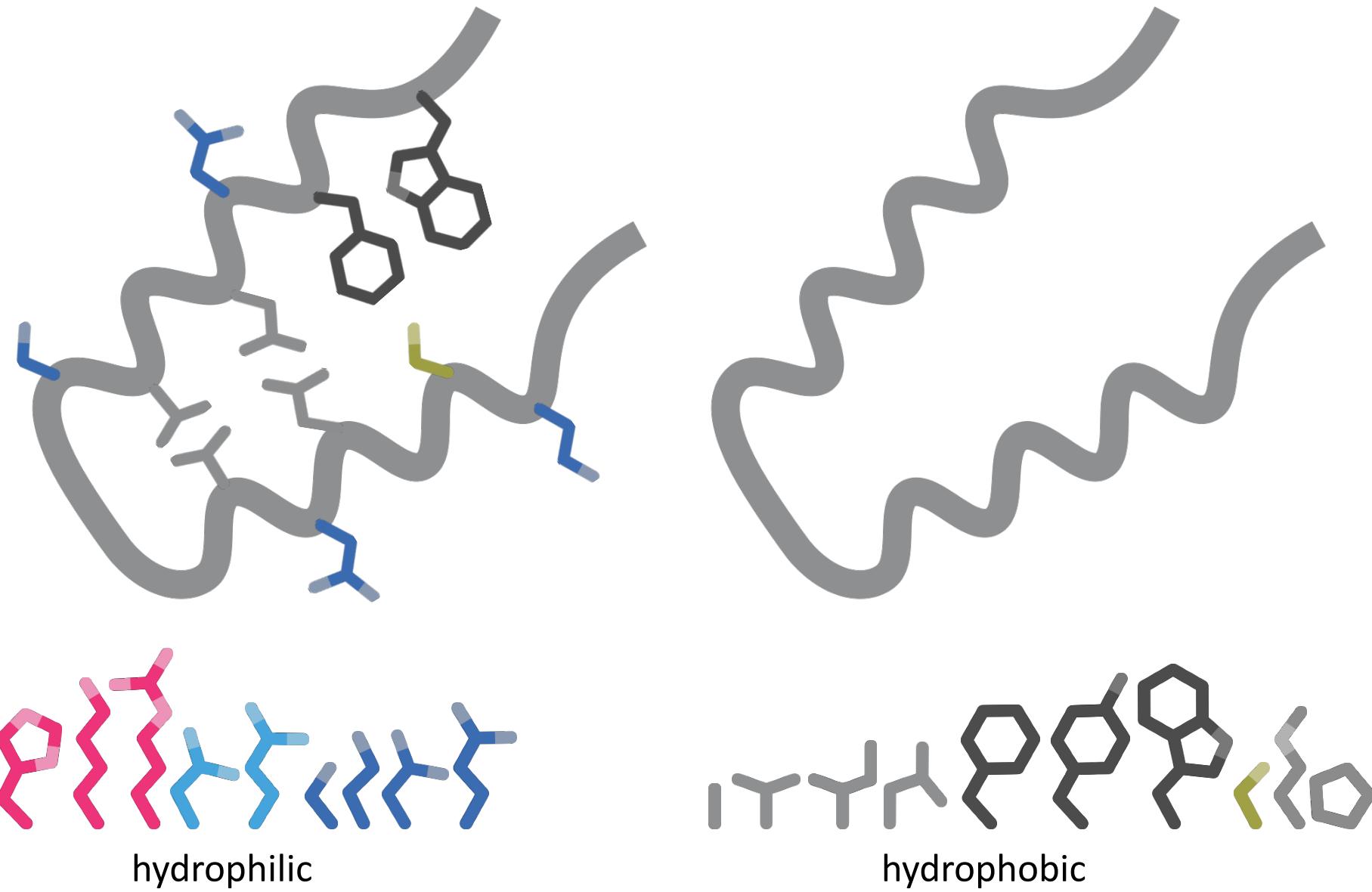
hydrophilic

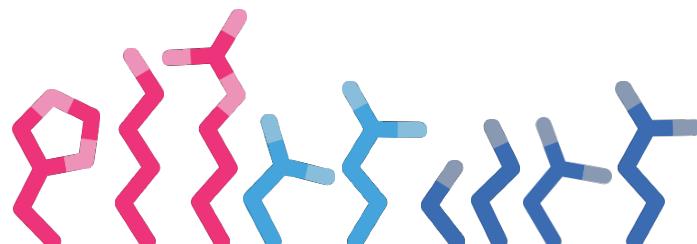
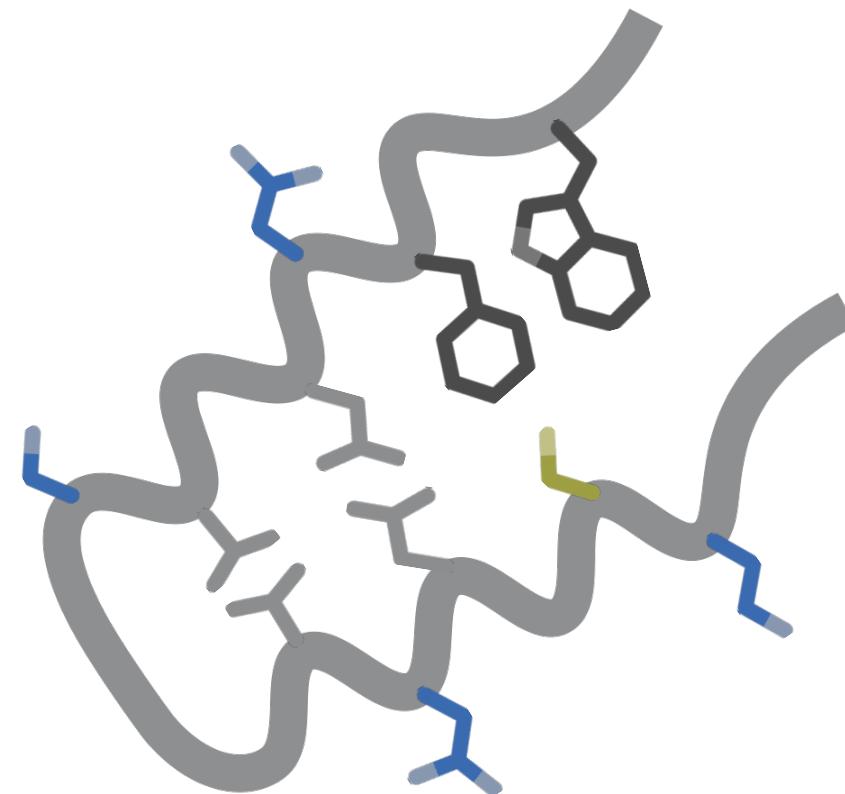


hydrophobic

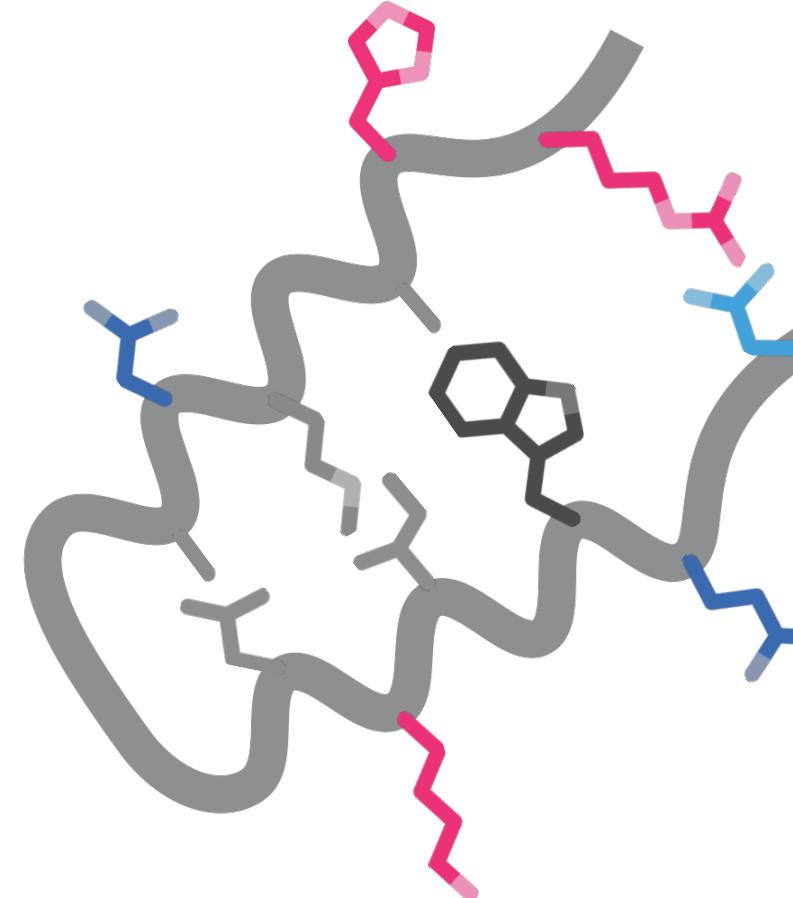




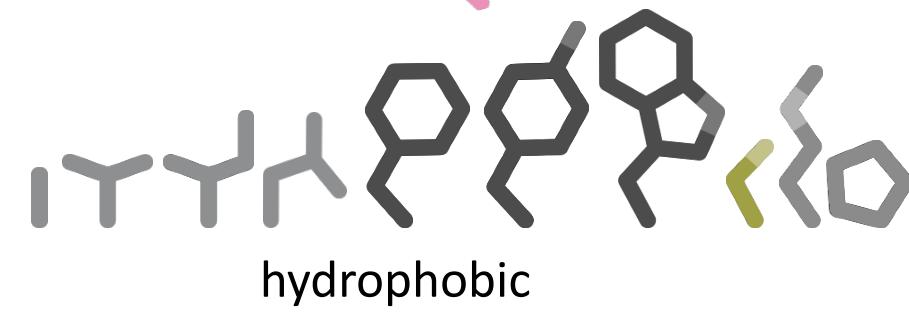
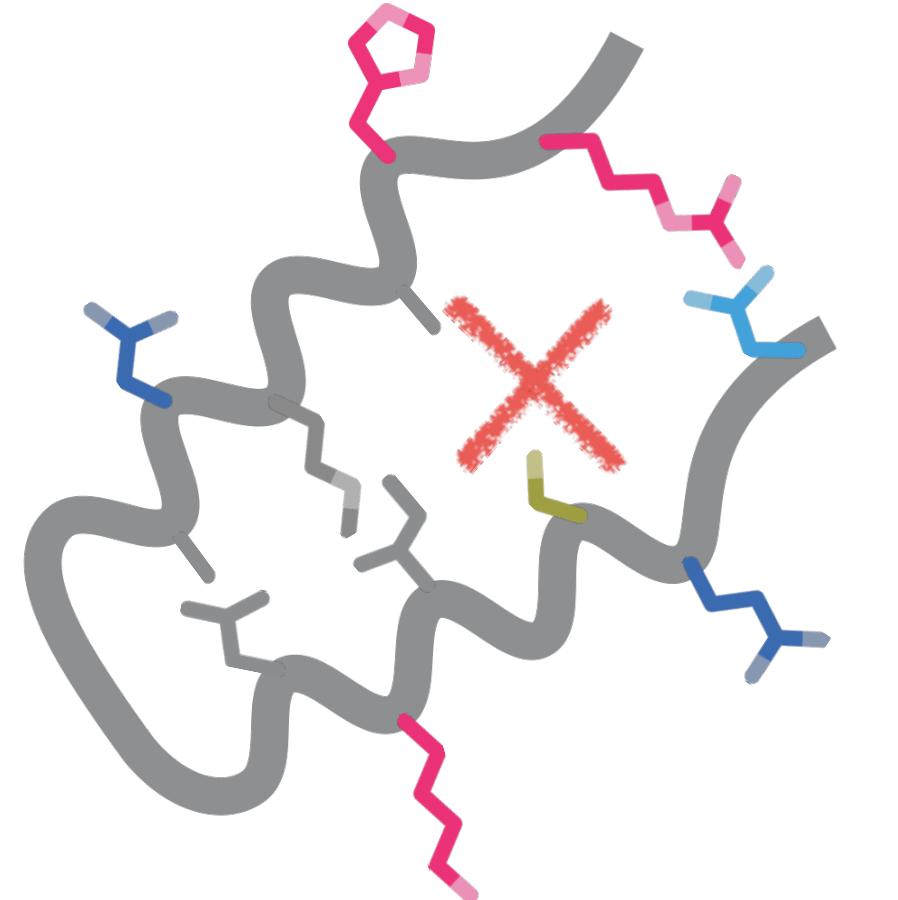
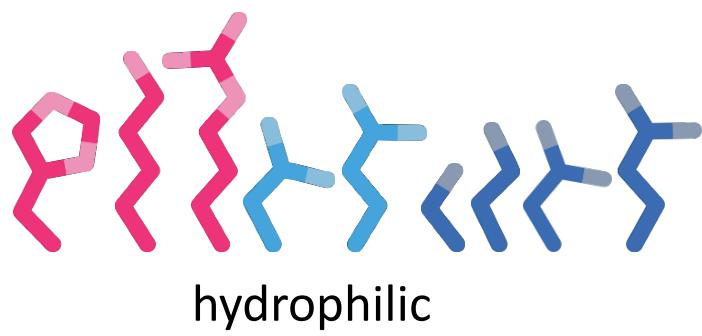
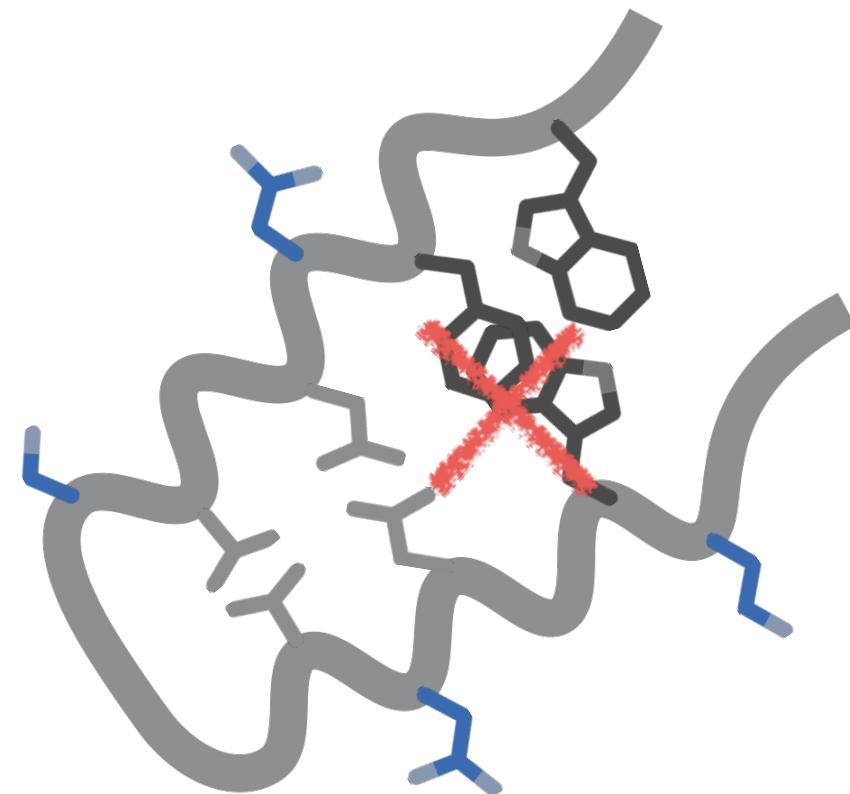


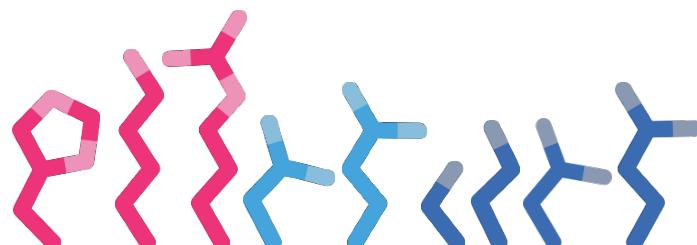
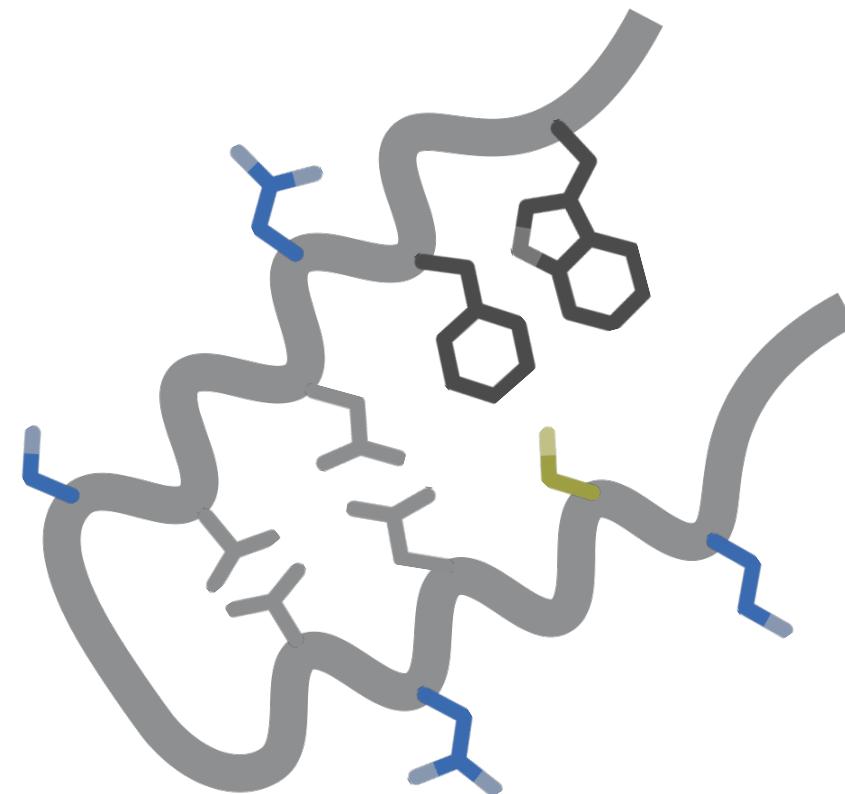


hydrophilic

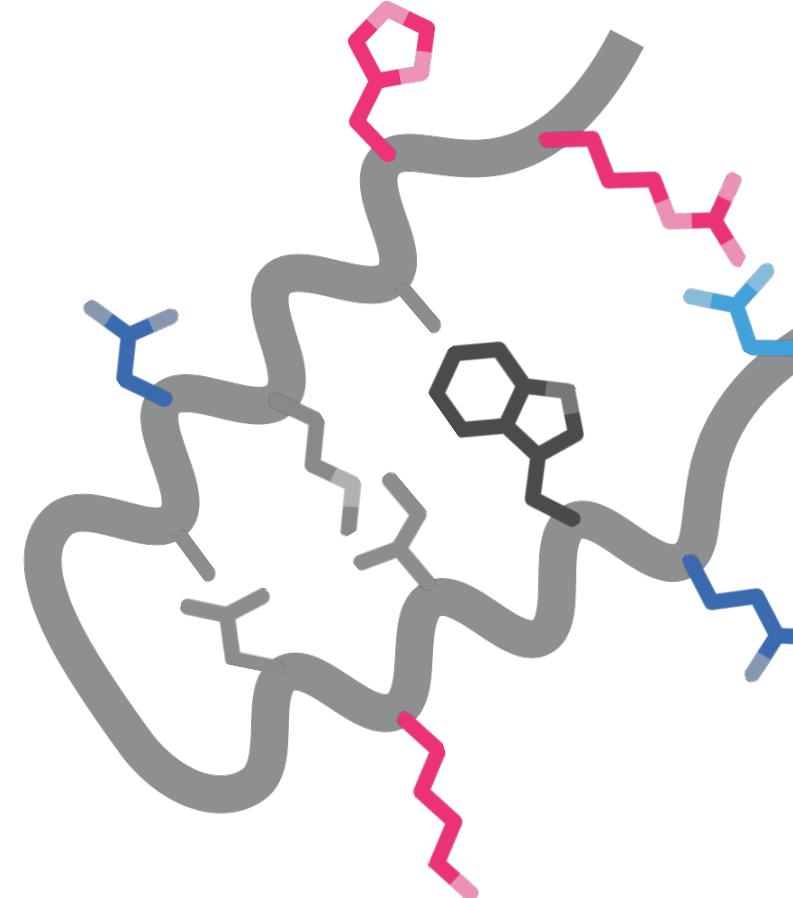


hydrophobic



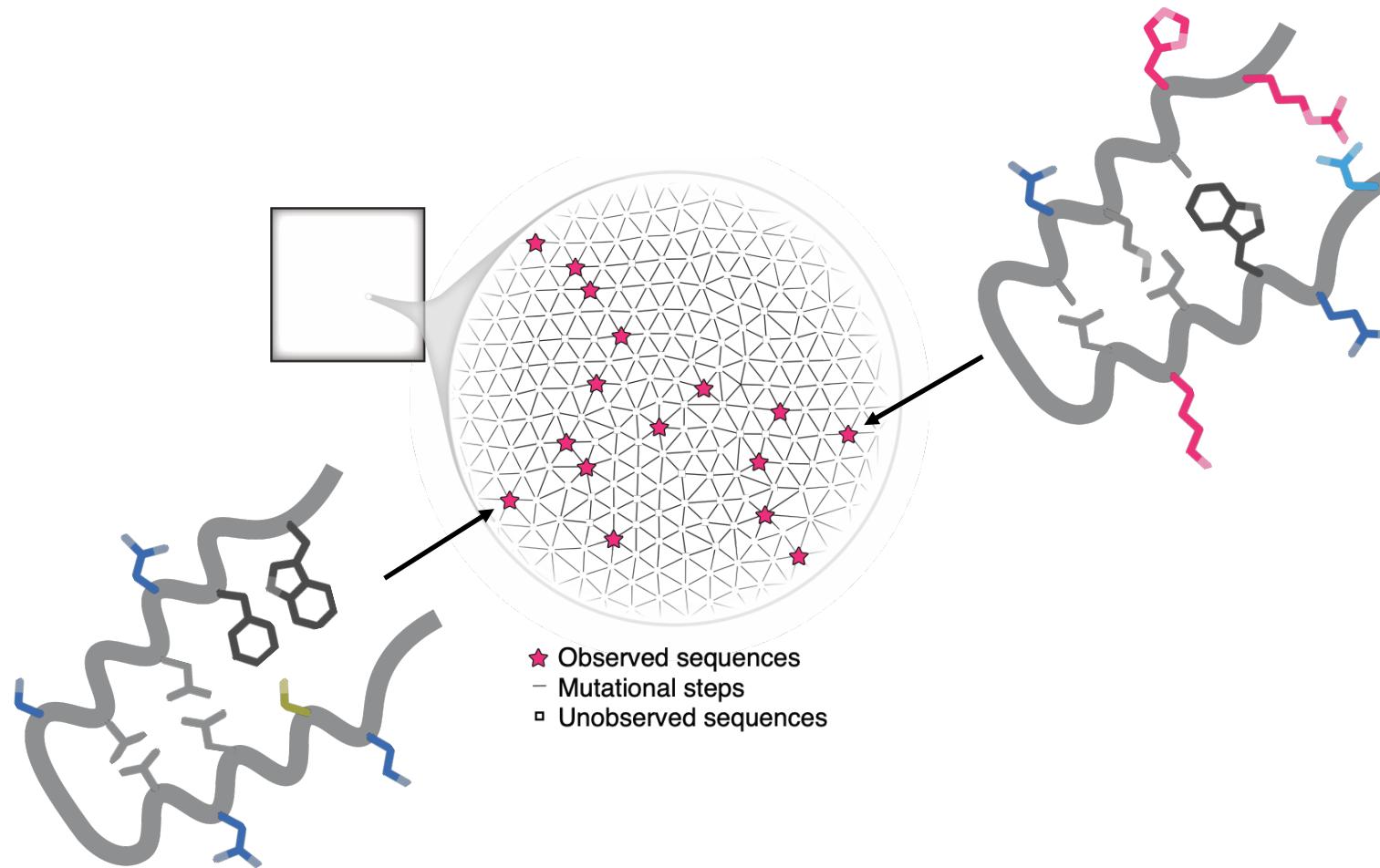


hydrophilic

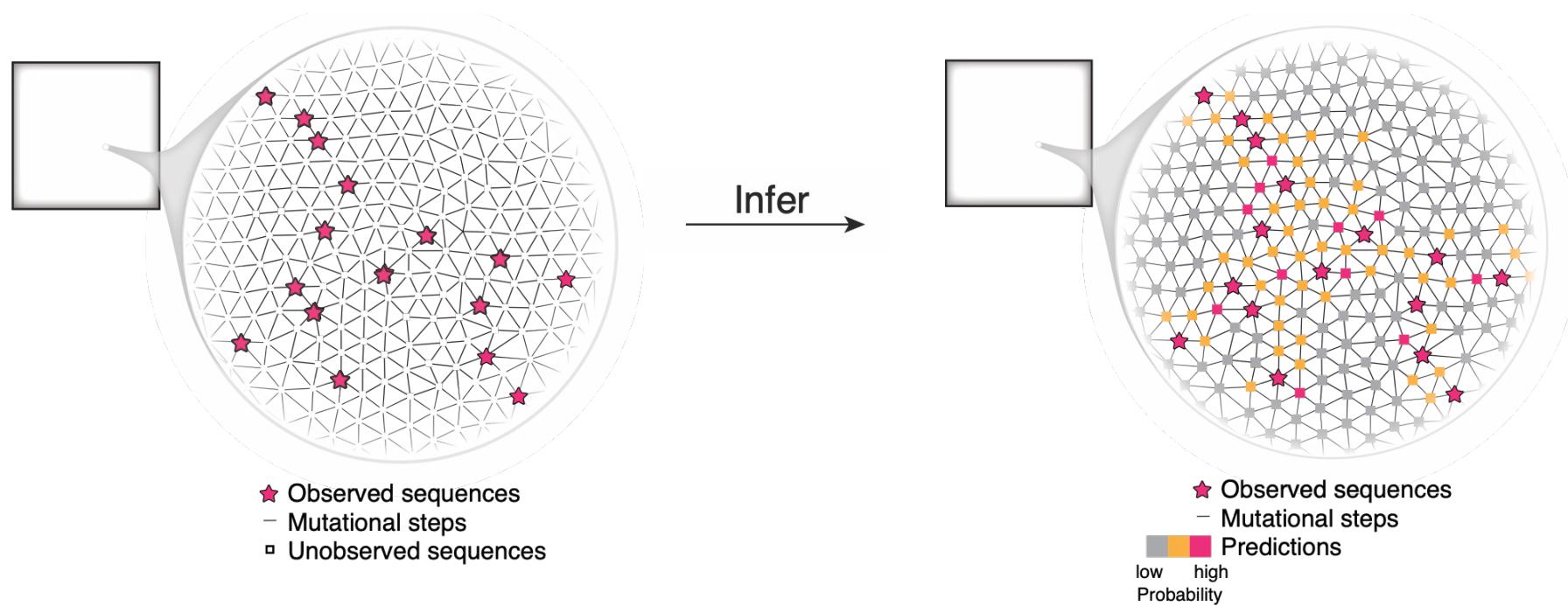


hydrophobic

# Evolution samples functional sequences



# We want a model that can learn the “patterns” of functional sequences



# The data: multiple sequence alignments

- ★ AQKLYLTHIDAEVEGD
- ★ ADRLYMTKIHHQFDGD
- ★ ADTLFITEVKQVFEGD
- ★ ADRLYMTKIHHTFDGD
- ★ ADKLYCTLIHNSFEGD
- ★ ADRLYMTKIHHEFEGD
- ★ ADRLYLTMIHQKFEAD
- ★ TDRLYITHIDETFEGD
- ★ ADRLYLTQIRNKFKGD
- ★ TSKMYITKIGQEFEGD
- ★ ADRLYMTKIHHEFEGD
- ★ ADRLYITHIHHHSFEGD
- ★ ADRLYMTKIHHEFEGD

- Structure and functional constraints is encoded in sequences
- Evolutionarily-related sequences can be aligned

# Many ways to model sequences

AQKLYLTHIDAEVEGD  
ADR<sup>LYMTKIHHQF</sup>DGD  
ADTLFITEV<sup>KQVF</sup>EGD  
ADR<sup>LYMTKIHHT</sup>F<sup>DGD</sup>  
AD<sup>KLYCTL</sup>I<sup>HNS</sup>FE<sup>GD</sup>  
ADR<sup>LYMTKIHHE</sup>FE<sup>GD</sup>  
ADR<sup>LYLTMIHQK</sup>FEAD  
TDR<sup>LYITHI</sup>DET<sup>FE</sup>GD  
ADR<sup>LYLTQIRNK</sup>FKGD  
TSK<sup>MYITKIGQE</sup>FE<sup>GD</sup>  
ADR<sup>LYMTKIHHE</sup>FE<sup>GD</sup>  
ADR<sup>LYITHIHHS</sup>FE<sup>GD</sup>  
ADR<sup>LYMTKIHHE</sup>FE<sup>GD</sup>

- Site-independent models
- Potts model (pairwise interactions)
- Higher order interactions

# Many ways to model sequences

AQKLYLTHIDAEVEGD  
ADR<sup>LYMTKIHHQFDGD</sup>  
ADTLFITEV<sup>KQVFEGD</sup>  
ADR<sup>LYMTKIHHTFDGD</sup>  
ADKLYCTL<sup>IHNSFEGD</sup>  
ADR<sup>LYMTKIHHFEGD</sup>  
ADR<sup>LYLTMIHQKF</sup>EAD  
TDR<sup>LYITHIDETFEGD</sup>  
ADR<sup>LYLTQIRNKFKGD</sup>  
TSKMYIT<sup>KIGQEFEVD</sup>  
ADR<sup>LYMTKIHHFEGD</sup>  
ADR<sup>LYITHIHHSFEGD</sup>  
ADR<sup>LYMTKIHHFEGD</sup>

- Site-independent models
- Potts model (pairwise interactions)
- Higher order interactions

# Site independence model

Multiple sequence alignment

AQ**K**LYLTHIDAEVEGD  
ADR**R**LYMTKIHHQFDGD  
ADT**L**FITEVKQVFEVD  
ADR**R**LYMTKIHHTFDGD  
AD**K**LYCTLIHNSFEGD  
ADR**R**LYMT**K**IHHEFEGD  
ADR**R**LYLTMIHQKFEAD  
TDR**R**LYITHIDETFEVD  
ADR**R**LYLTQIRNKFKGD  
TS**K**MYITKIGQEFEVD  
ADR**R**LYMT**K**IHHEFEGD  
ADR**R**LYITHIHHSFEGD  
ADR**R**LYMT**K**IHHEFEGD

Position  $i$

Position specific amino acid frequencies

$$h_1 \quad h_2 \quad \dots \quad h_i \quad \dots \quad h_L$$

$$h_i = \begin{array}{c} | \\ \text{A} \\ \text{C} \\ \text{D} \\ \text{E} \\ \text{F} \\ \text{G} \\ \text{H} \\ \text{I} \\ \text{K} \\ \text{L} \\ \text{M} \\ \text{N} \\ \text{P} \\ \text{Q} \\ \text{R} \\ \text{S} \\ \text{T} \\ \text{V} \\ \text{W} \\ \text{Y} \end{array}$$

Score sequences

$$\boldsymbol{x} = \text{ADR} \text{LYMT} \text{K} \text{IHHQFDGT}$$

Product of site factors

$$P(\boldsymbol{x}) = p_1(\text{A})p_2(\text{D})p_3(\text{R}) \dots p_L(\text{T})$$

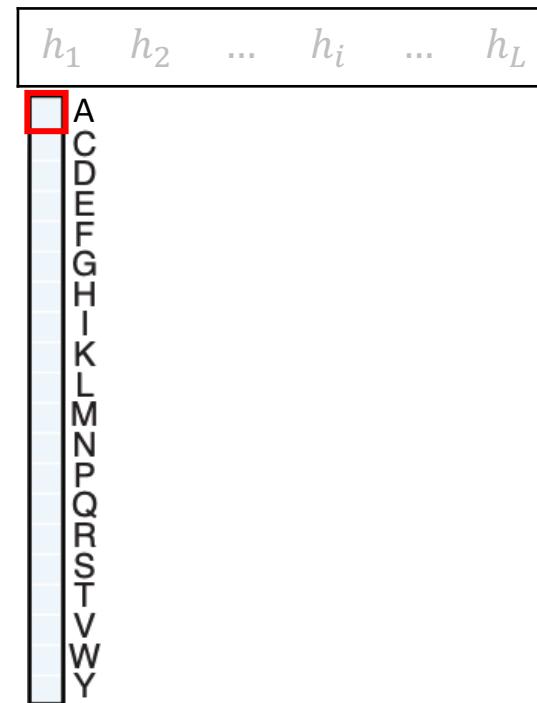
# Site independence model

Multiple sequence alignment

AQ**K**LYLTHIDAEVEGD  
ADR**R**LYMTKIHHQFDGD  
ADT**L**FITEVKQVFEGD  
ADR**R**LYMTKIHHTFDGD  
AD**K**LYCTLIHNSFEGD  
ADR**R**LYMT**K**IHHEFEGD  
ADR**R**LYLTMIHQKFEAD  
TDR**R**LYITHIDETFEGD  
ADR**R**LYLTQIRNKFKGD  
TS**K**MYITKIGQEFEVD  
ADR**R**LYMT**K**IHHEFEGD  
ADR**R**LYITHIHHSFEGD  
ADR**R**LYMT**K**IHHEFEGD

Position  $i$

Position specific amino acid frequencies



Score sequences

$$\boldsymbol{x} = \boxed{A} \text{DRLYMT} \boxed{K} \text{IHHQFDGT}$$

Product of site factors

$$P(\boldsymbol{x}) = p_1(A)p_2(D)p_3(R) \dots p_L(T)$$

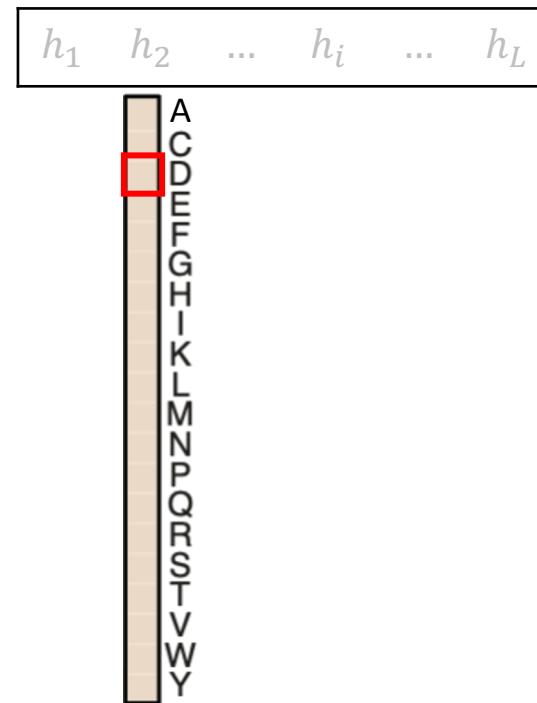
# Site independence model

Multiple sequence alignment

AQ**K**LYLTHIDAEVEGD  
ADR**R**LYMTKIHHQFDGD  
ADT**L**FITEVK**Q**VFEVD  
ADR**R**LYMTKI**H**HHTFDGD  
AD**K**LY**C**TLLIHNSFEGD  
ADR**R**LYMT**K**I**H**HEFEGD  
ADR**R**LYLTMI**H**QKFEAD  
TDR**R**LYITHIDETFEGD  
ADR**R**LYLT**Q**IRNKFKGD  
TS**K**MYIT**K**IGQEFEVD  
ADR**R**LYMT**K**I**H**HEFEGD  
ADR**R**LYITHI**H**HHSFEGD  
ADR**R**LYMT**K**I**H**HEFEGD

Position  $i$

Position specific amino acid frequencies



Score sequences

$$\boldsymbol{x} = A\boxed{D}R\textcolor{blue}{L}Y\textcolor{red}{M}\textcolor{blue}{T}\textcolor{red}{K}\textcolor{blue}{I}\textcolor{red}{H}\textcolor{blue}{H}\textcolor{red}{Q}\textcolor{blue}{F}\textcolor{red}{D}\textcolor{blue}{G}\textcolor{red}{T}$$

Product of site factors

$$P(\boldsymbol{x}) = p_1(A)p_2(D)p_3(R) \dots p_L(T)$$

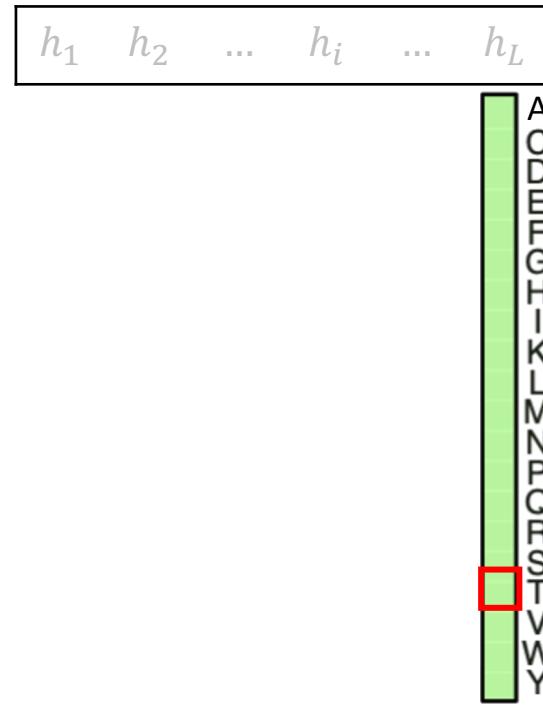
# Site independence model

Multiple sequence alignment

AQ**K**LYLTHIDAEVEGD  
ADR**R**LYMT**KI**HHQFDGD  
ADT**L**FITEV**KQ**VFEGD  
ADR**R**LYMT**KI**HHTFDGD  
AD**K**LY**C**TLLIHNSFEGD  
ADR**R**LYMT**KI**HHEFEGD  
ADR**R**LYLTMI**HQ**KFEAD  
TDR**R**LYITHIDETFEGD  
ADR**R**LYLT**QI**RNKFKGD  
TSK**M**YIT**KI**GQEFEVD  
ADR**R**LYMT**KI**HHEFEGD  
ADR**R**LYITHI**HHS**FEGD  
ADR**R**LYMT**KI**HHEFEGD

Position  $i$

Position specific amino acid frequencies



Score sequences

$$\boldsymbol{x} = \text{ADR} \text{LYM} \text{T} \text{KI} \text{HHQ} \text{FDG} \text{T}$$

Product of site factors

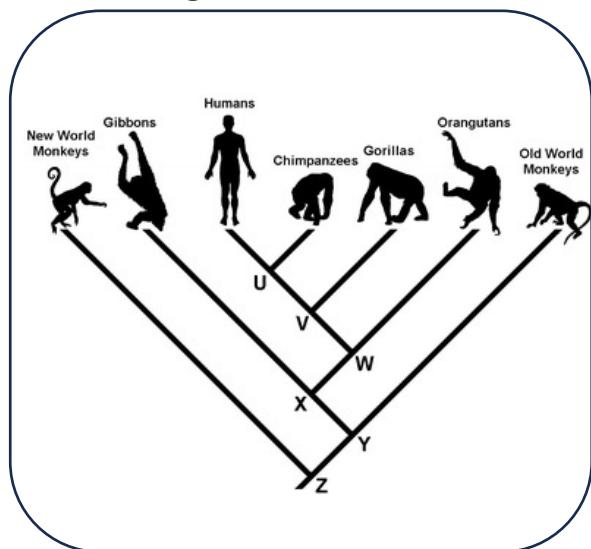
$$\begin{aligned} P(\boldsymbol{x}) &= p_1(A)p_2(D)p_3(R) \dots p_L(T) \\ &= p_1(x_1)p_2(x_2) \dots p_L(x_L) \end{aligned}$$

or equivalently

$$\begin{aligned} P(\boldsymbol{x}) &= \frac{1}{Z} \exp\left(\sum_i h_i(x_i)\right) \\ h_i &= \log(p_i(x_i)) \end{aligned}$$

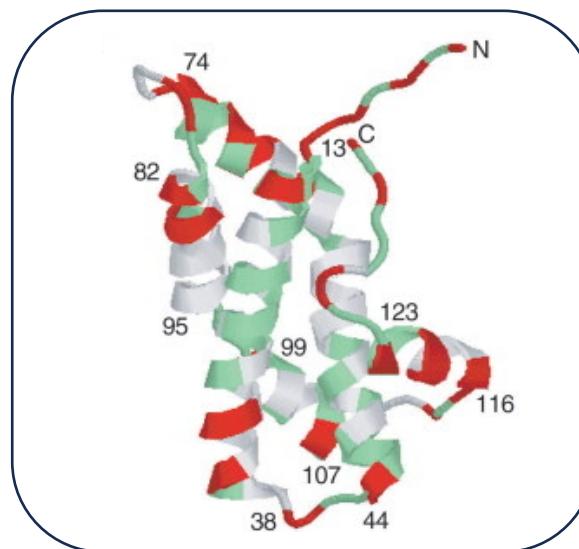
# Site independence model

Phylogenetic inference



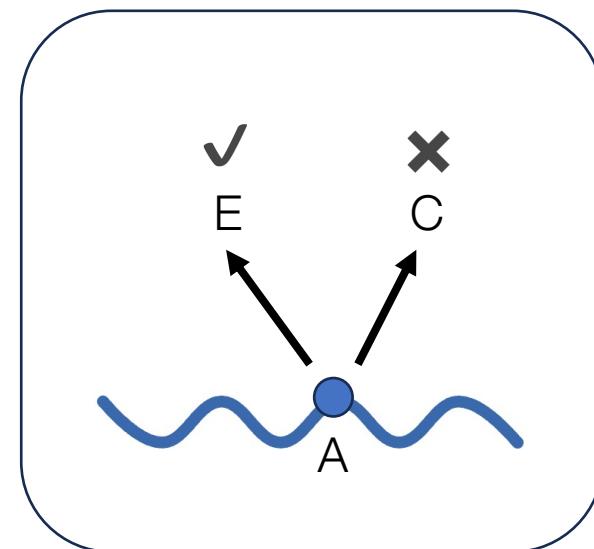
Ronquist *et al* 2012 *Syst Biol*  
Suchard *et al* 2018 *Virus Evol*  
Minh *et al* 2020 *Mol Biol Evol*

Substitution rate estimates



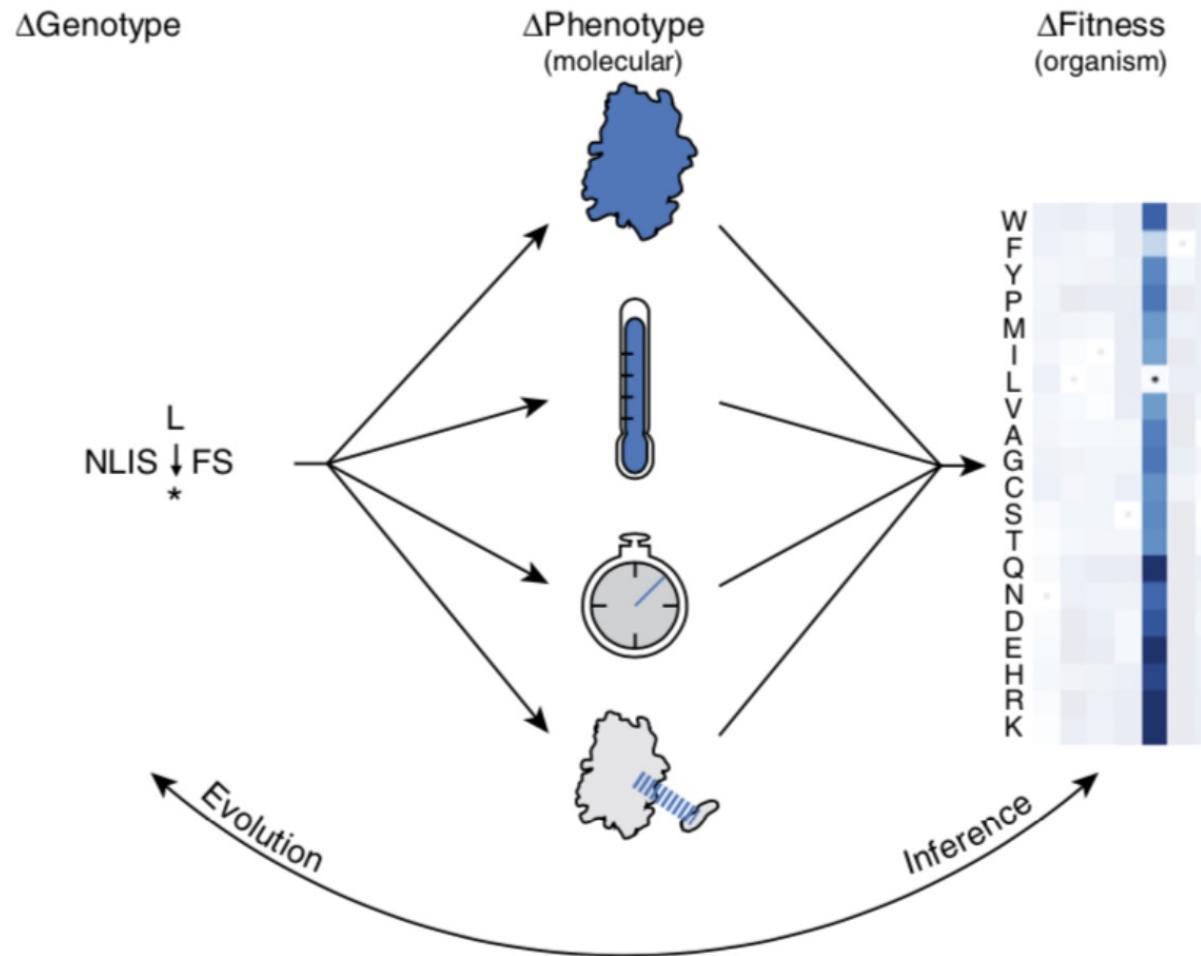
Pond *et al* 2020 *Mol Biol Evol*  
Jones, Youssef *et al* 2020 *Syst Biol*  
Youssef *et al* 2020 *Mol Biol Evol*

Mutation effect prediction

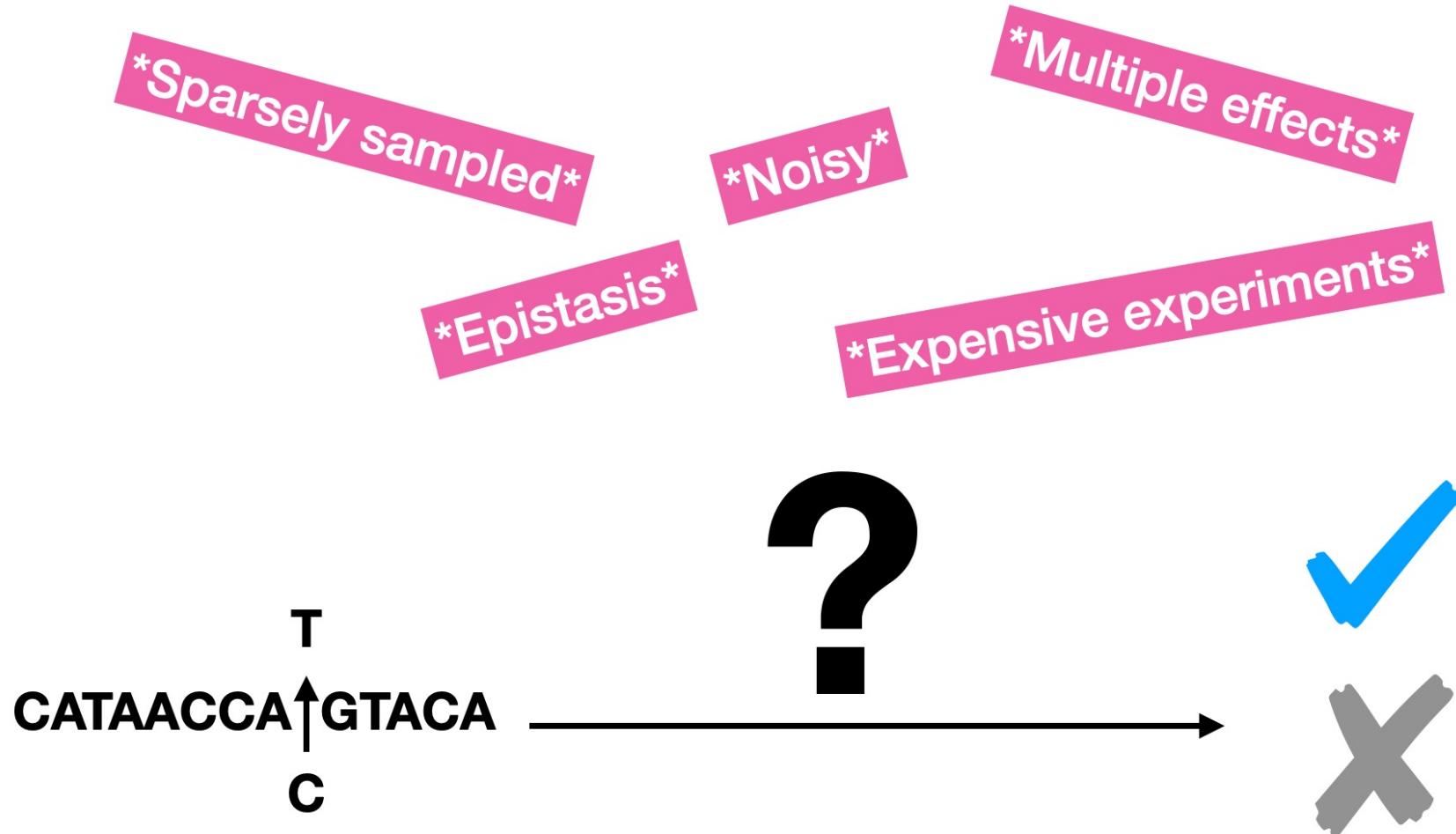


Hopf *et al* 2017 *Nature Biotech*  
Riesselman *et al* 2018 *Nature Meth*

# Mutation effect prediction: Mutations alter biological function



# Mutation effect prediction is hard!



# Site independence model

**Most common name**

**Mohammed**

---

First name

**James**

---

Middle name

**Wang**

---

Last name

Epistasis: Sequence context matters!

# Many ways to model sequences

AQKLYLTHIDAEVEGD  
ADR<sup>LYMTKIHHQFDGD</sup>  
ADTLFITEV<sup>KQVFEGD</sup>  
ADR<sup>LYMTKIHHTFDGD</sup>  
ADKLYCTL<sup>IHNSFEGD</sup>  
ADR<sup>LYMTKIHHFEGD</sup>  
ADR<sup>LYLTMIHQKF</sup>EAD  
TDR<sup>LYITHIDETFEGD</sup>  
ADR<sup>LYLTQIRNKFKGD</sup>  
TSKMYIT<sup>KIGQEFEVD</sup>  
ADR<sup>LYMTKIHHFEGD</sup>  
ADR<sup>LYITHIHHSFEGD</sup>  
ADR<sup>LYMTKIHHFEGD</sup>

- Site-independent models
- Potts model (pairwise interactions)
- Higher order interactions

# Pairwise interactions model

Multiple sequence alignment

A multiple sequence alignment showing 12 rows of amino acid sequences. The first row starts with 'AQ' and ends with 'GD'. The second row starts with 'ADR' and ends with 'GD'. The third row starts with 'ADT' and ends with 'GD'. The fourth row starts with 'ADR' and ends with 'GD'. The fifth row starts with 'ADK' and ends with 'GD'. The sixth row starts with 'ADR' and ends with 'GD'. The seventh row starts with 'ADR' and ends with 'GD'. The eighth row starts with 'TDR' and ends with 'GD'. The ninth row starts with 'ADR' and ends with 'GD'. The tenth row starts with 'TSK' and ends with 'GD'. The eleventh row starts with 'ADR' and ends with 'GD'. The twelfth row starts with 'ADR' and ends with 'GD'. Two specific positions are highlighted: position  $i$  (the 4th column) and position  $j$  (the 10th column). These highlighted positions are enclosed in black rectangular boxes.

Position  $i$

Position specific amino acid frequencies

A horizontal vector representing position-specific amino acid frequencies. It consists of a series of boxes containing labels:  $h_1$ ,  $h_2$ , ...,  $h_i$ , ...,  $h_L$ . The labels are separated by ellipses indicating intermediate positions.

Position  $j$

# Pairwise interactions model

Multiple sequence alignment

AQKLYLTHIDAEVEGD  
ADR~~LYM~~T~~KI~~HHQFDGD  
ADTLFITEV~~K~~QVFEGD  
ADR~~LYM~~T~~KI~~HHTFDGD  
AD~~K~~LYCTL~~I~~HNSFE~~G~~D  
ADR~~LYM~~T~~KI~~HHEFE~~G~~D  
ADR~~LYL~~TMIHQKF~~EAD~~  
TDR~~LYI~~THIDETFE~~G~~D  
ADR~~LYL~~TQIRNKF~~K~~GD  
TSKMYIT~~K~~IGQE~~F~~GD  
ADR~~LYM~~T~~KI~~HHEFE~~G~~D  
ADR~~LYI~~THIHH~~S~~FE~~G~~D  
ADR~~LYM~~T~~KI~~HHEFE~~G~~D

Position *i*      Position *j*

Position specific amino acid frequencies

$$h_1 \quad h_2 \quad \dots \quad h_i \quad \dots \quad h_L$$

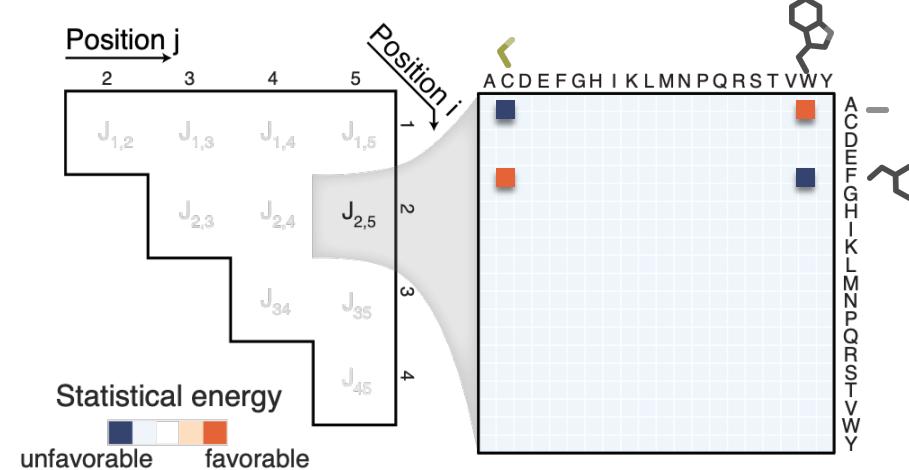
Score sequences

ADR~~LYM~~T~~KI~~HHQFDGD

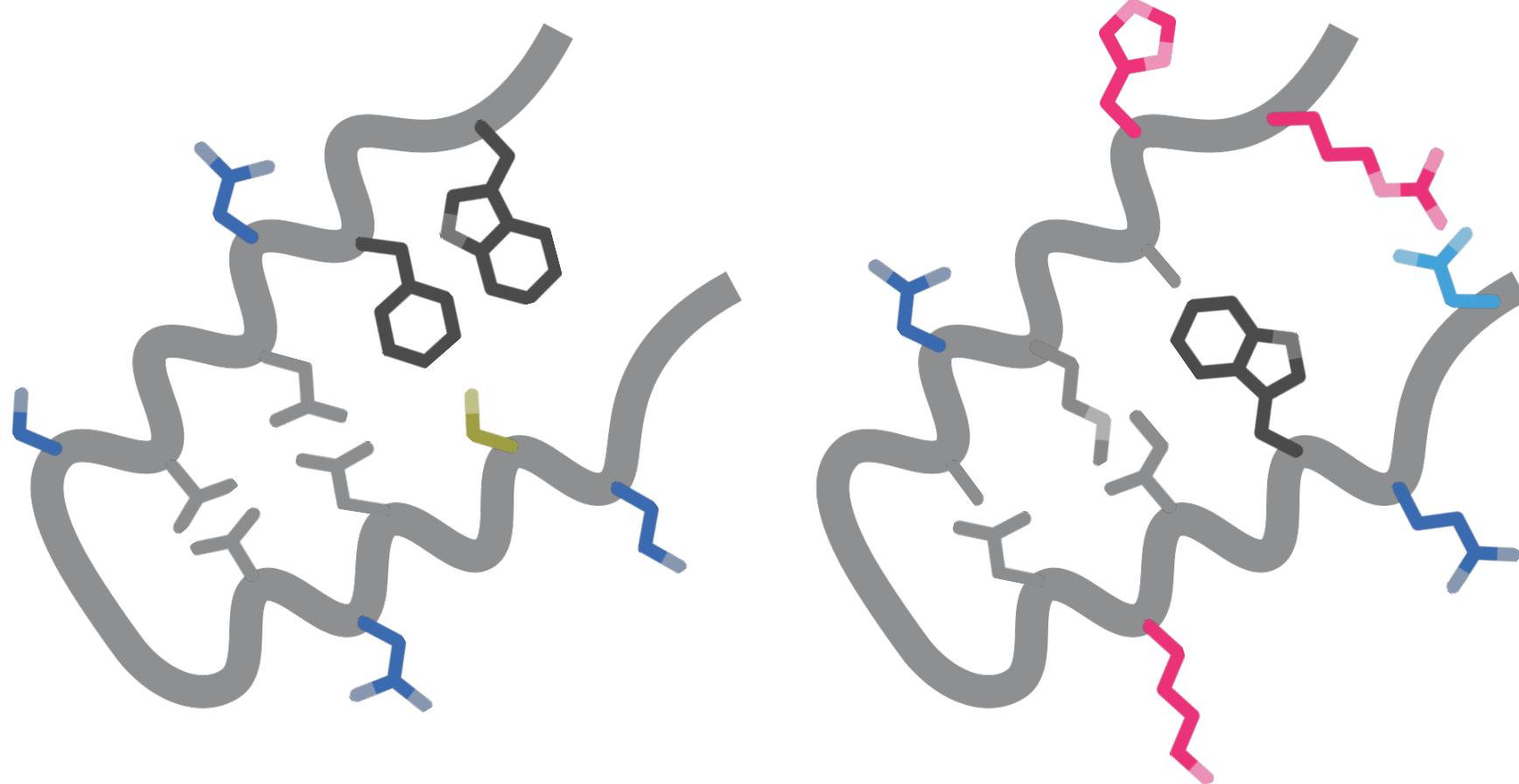
Product of **site** and **pair** factors

$$P(x) = \frac{1}{Z} \exp \left( \sum_i h_i(x_i) + \sum_{i < j} j_{ij}(x_i, x_j) \right)$$

Pairwise frequencies

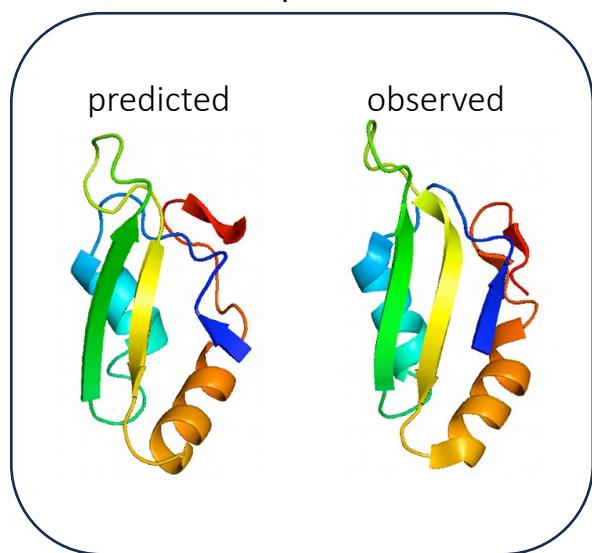


# Co-evolving sites (correlated mutations) are likely 3D contacts



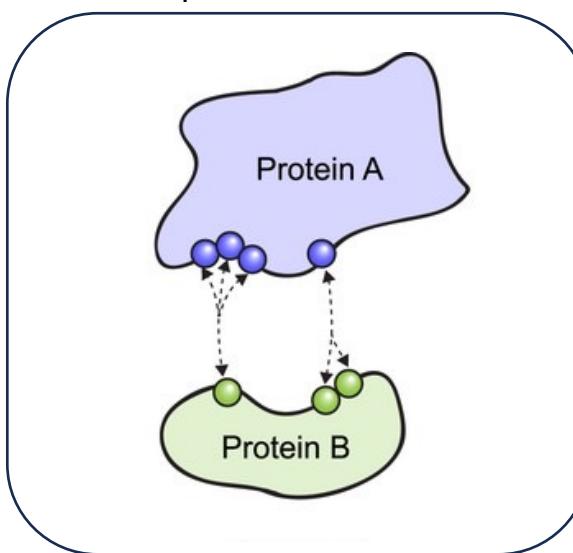
# Pairwise interactions model

Structure prediction



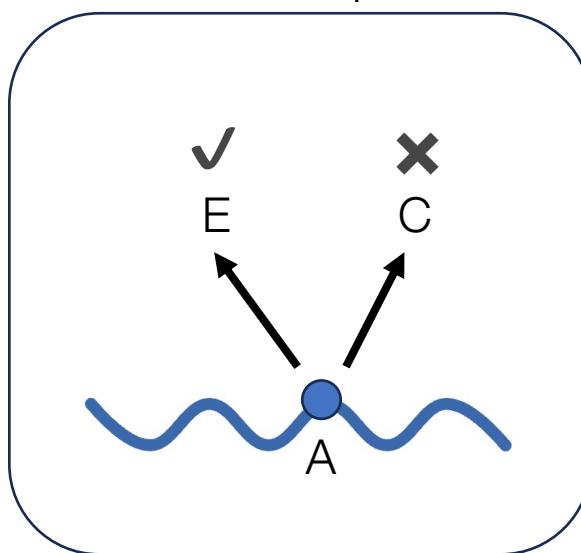
Hopf *et al* 2012 *Cell*  
Marks *et al* 2012 *Nature Biotech*

Protein-protein interactions



Green *et al* 2021 *Nature Comm*  
Hopf *et al* 2014 *eLife*

Mutation effect prediction



Hopf *et al* 2017 *Nature Biotech*  
Riesselman *et al* 2018 *Nature Meth*

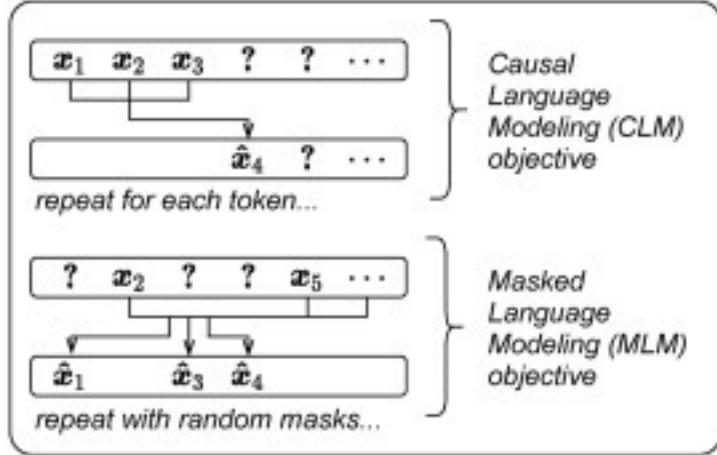
# Many ways to model sequences

AQKLYLTHIDAEVEGD  
ADR<sup>LYMTKIHHQFDGD</sup>  
ADTLFITEV<sup>KQVFEGD</sup>  
ADR<sup>LYMTKIHHHTFDGD</sup>  
ADKLYCTL<sup>IHNSFEGD</sup>  
ADR<sup>LYMTKIHHFEGD</sup>  
ADR<sup>LYLTMIHQKF</sup>EAD  
TDR<sup>LYITHIDETFEGD</sup>  
ADR<sup>LYLTQIRNKFKGD</sup>  
TSKMYIT<sup>KIGQEFEVD</sup>  
ADR<sup>LYMTKIHHFEGD</sup>  
ADR<sup>LYITHIHHSFEGD</sup>  
ADR<sup>LYMTKIHHFEGD</sup>

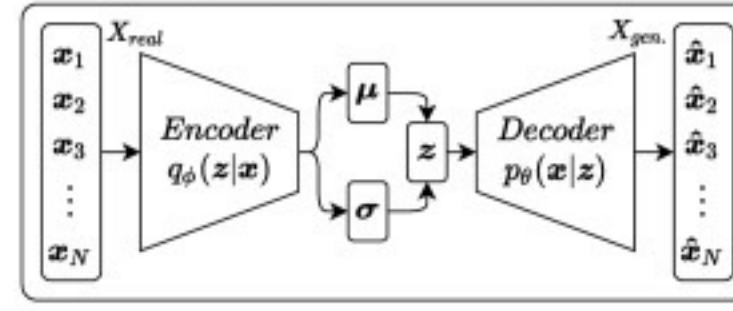
- Site-independent models
- Potts model (pairwise interactions)
- Higher order interactions

# Many ways to model higher order interactions

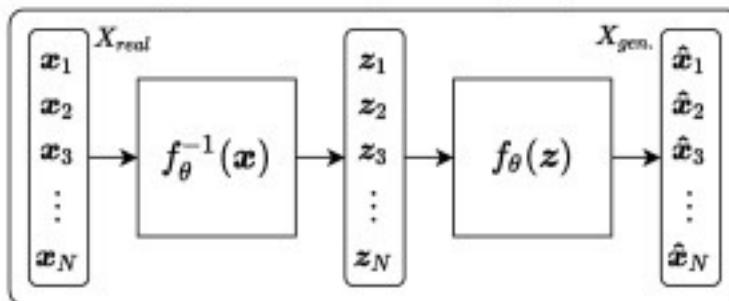
Autoregressive Models



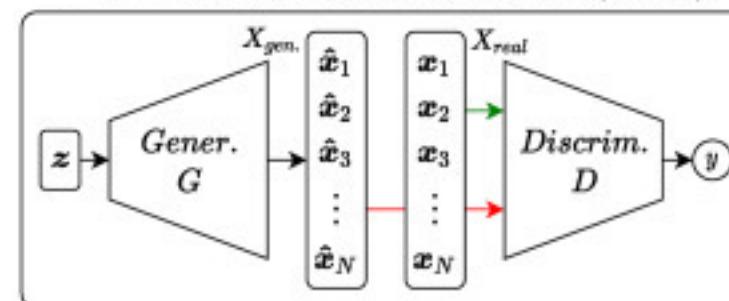
Variational Autoencoders (VAEs)



Normalising Flows (NFs)



Generative Adversarial Networks (GANs)



# Protein language models don't need multiple sequence alignments

# Challenges

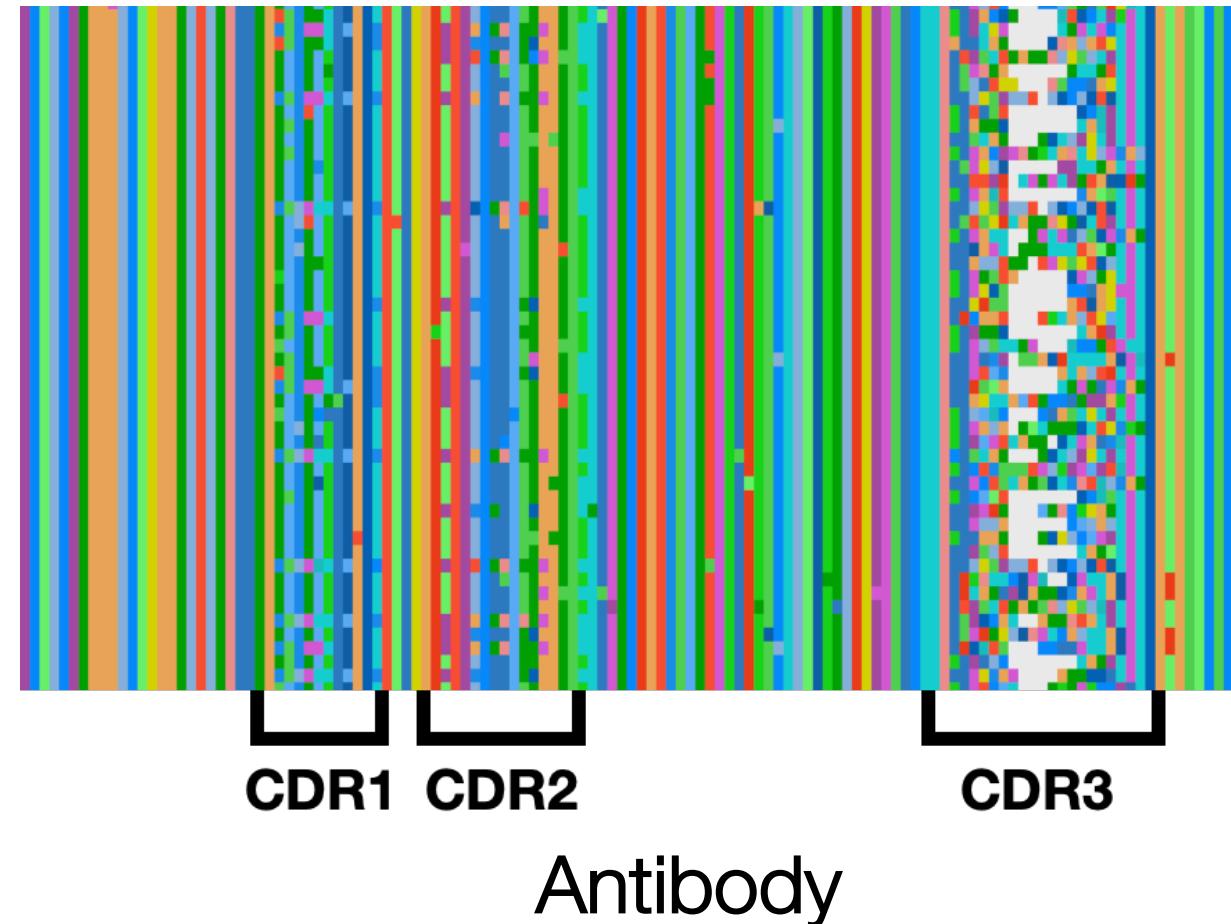
# Gap characters

# Alignment uncertainty

# Insertions and deletions

AQKLYLTHIDAEVEGD  
ADRLYMTKIHHQFDGD  
ADTLFITEVKQVFEGD  
--RLYMTKIHHTFD--  
ADKLYCTLIH-SFEGD  
ADRLY-TKIHHEFEGD  
--RLYLTMIHQKFEAD  
TDRLYITHIDETFEGD  
ADRLYLTQIRNKFKGD  
TSKMYITKIGQEFEGD  
--RLYMTKIH-EFEGD  
ADRLYITHIH-SFEGD  
--RLYMTKIH-EFEGD

# Some proteins like antibodies are just hard to align



If there is no alignment...we aren't limited to  
**one** protein family that aligns together.

What if we built a model of **all** proteins?

# Rise of large language models for proteins

Article | Published: 21 Oct 2019

## ProGen: Language Modeling for Protein Generation

Unified ration based deep re

Ali Madani<sup>1</sup> Bryan

Ethan C. Alley, Grigory K

Nature Methods 16, 1315–1322 (2019) | Cite thi

ZEMING LIN  , HALIL AKIN  , ROSHAN RAO  , BRIAN HIE  , ZHONGKAI ZHU, WENTING LU, NIKITA SMETANIN, ROBERT VERKUIL  , ORI KABELI  , [...], AND

### Embeddings from protein language models predict conservation and variant effects

Céline Marquet<sup>1,2</sup>  · Michael Heinzinger<sup>1,2</sup> · Tobias Olenyi<sup>1</sup>,  
Michael Bernhofer<sup>1,2</sup> · Dmitrii Nechaev<sup>1,2</sup> · Burkhard Rost<sup>1,3</sup>

Received: 1 June 2021 / Accepted: 6 December 2021 / Published online: 30 Dec 2021  
© The Author(s) 2021

### from scaling unsupervised learning million protein sequences

Alexander Rives   , Joshua Meier, Tom Sercu  ,  , and Rob Fergus [Authors Info & Affiliations](#)

Edited by David T. Jones, University College London, London, United Kingdom, and accepted by Editorial Board Member, December 16, 2020 (received for review August 6, 2020)

April 5, 2021 | 118 (15) e2016239118 | <https://doi.org/10.1073/pnas.2016239118>

## ProGen2: Exploring the Boundaries of Protein Language Models

Erik Nijkamp\*  
Salesforce Research

Jeffrey Ruffolo\*  
Johns Hopkins University

Eli N. Weinstein  
Columbia University

Nikhil Naik  
Salesforce Research

Ali Madani  
Salesforce Research

# Language modeling training with self-supervision (no labels) and without alignments

**Autoregressive** (predict next token):

Natural language: The brown fox jumped over the \_\_\_\_\_

Protein Sequence: GYAMLIVEQGA\_\_\_\_\_

**Masked** (predict masked tokens):

Natural language: The brown \_\_\_\_\_ jumped over the \_\_\_\_\_ dog

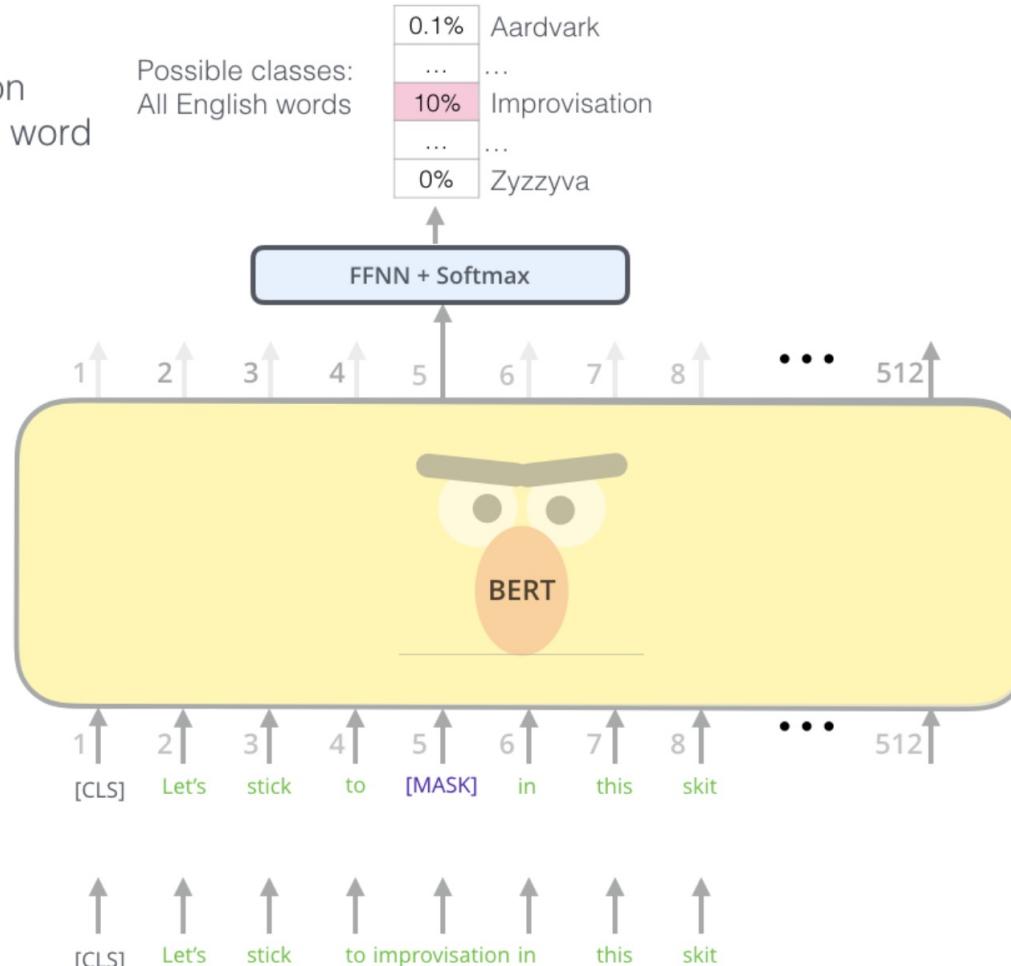
Protein Sequence: GYAM\_IVEEQGAL\_SKGV\_AIT

# Masked model (pre-)training from BERT

Use the output of the masked word's position to predict the masked word

Randomly mask 15% of tokens

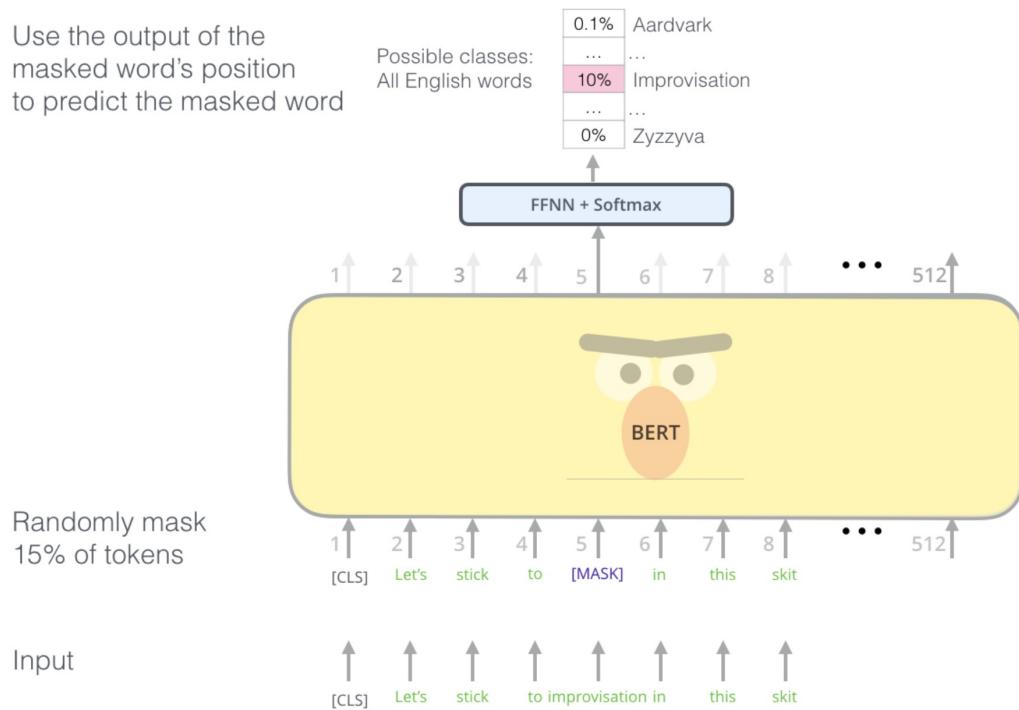
Input



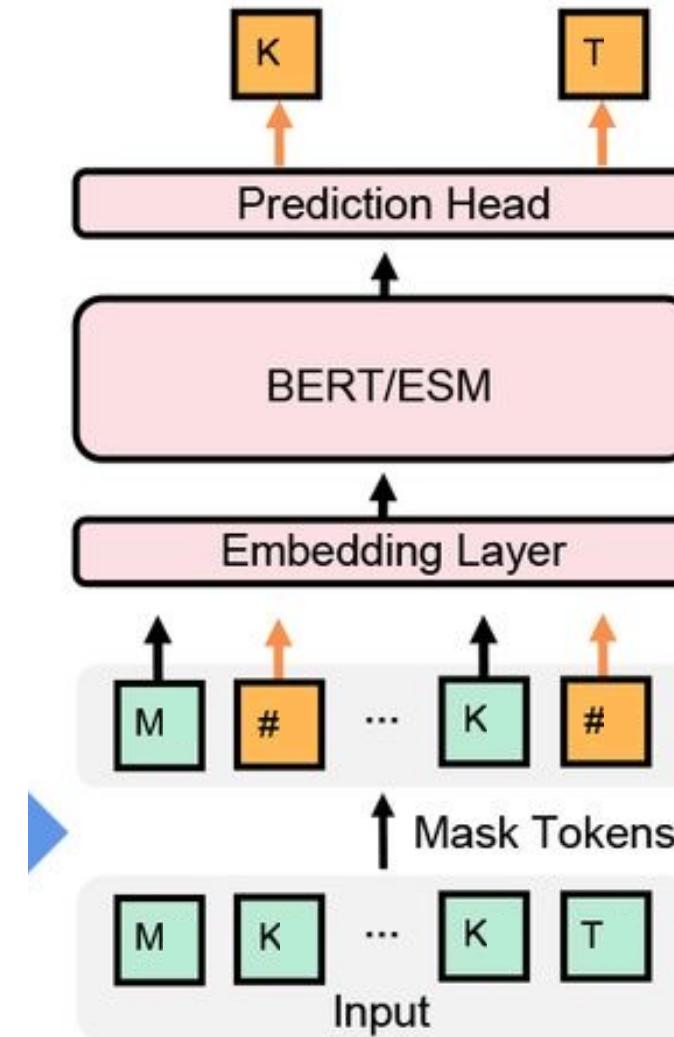
BERT's clever language modeling task masks 15% of words in the input and asks the model to predict the missing word.

# ESM uses the BERT framework

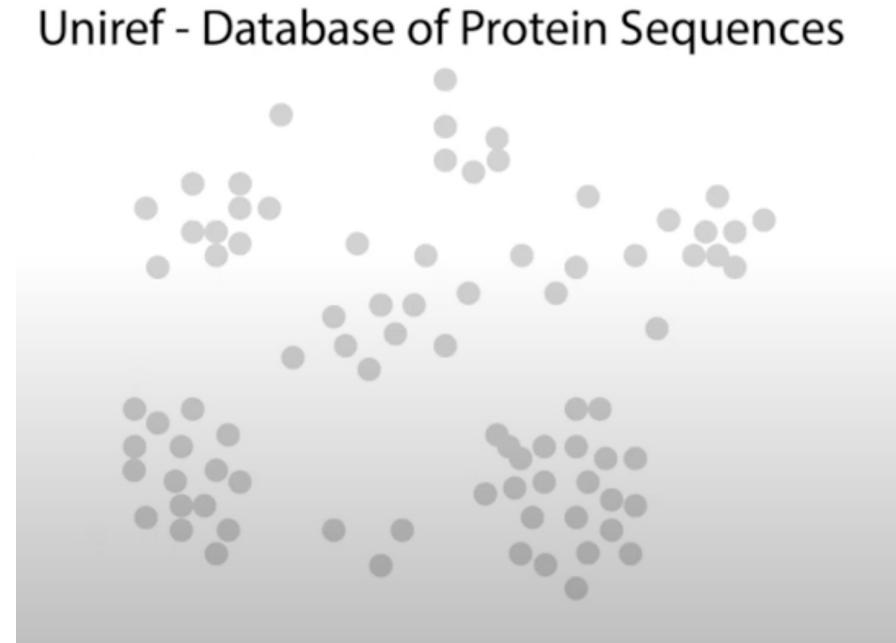
Use the output of the masked word's position to predict the masked word



BERT's clever language modeling task masks 15% of words in the input and asks the model to predict the missing word.

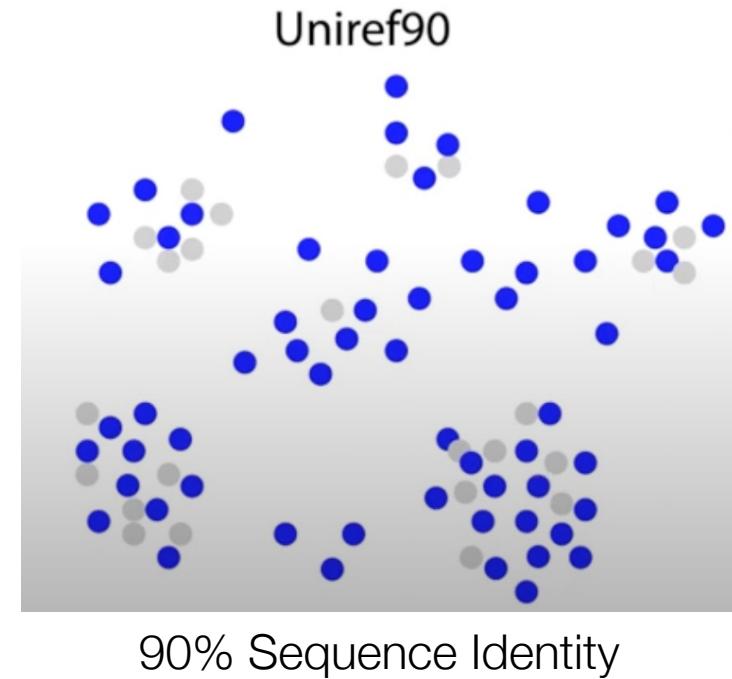
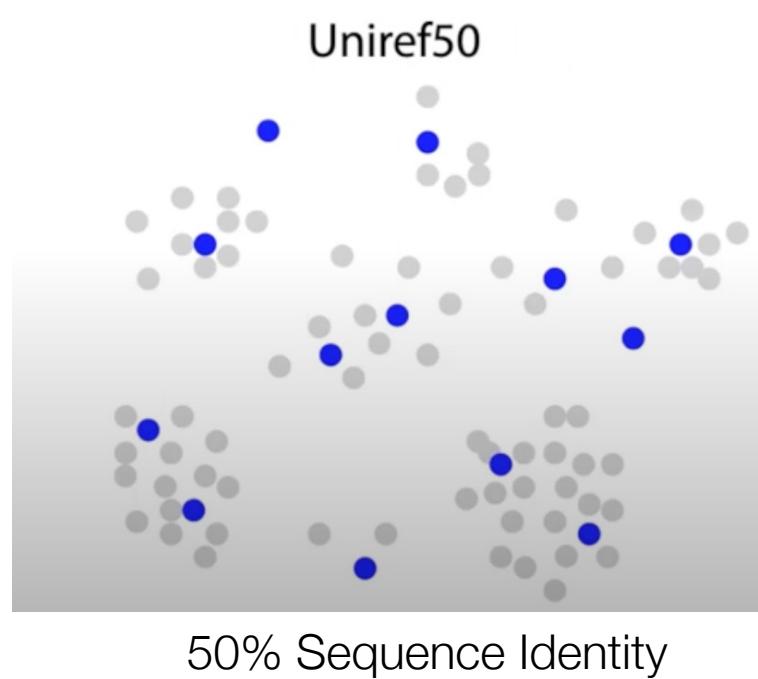


# Most protein language models are trained on UniRef



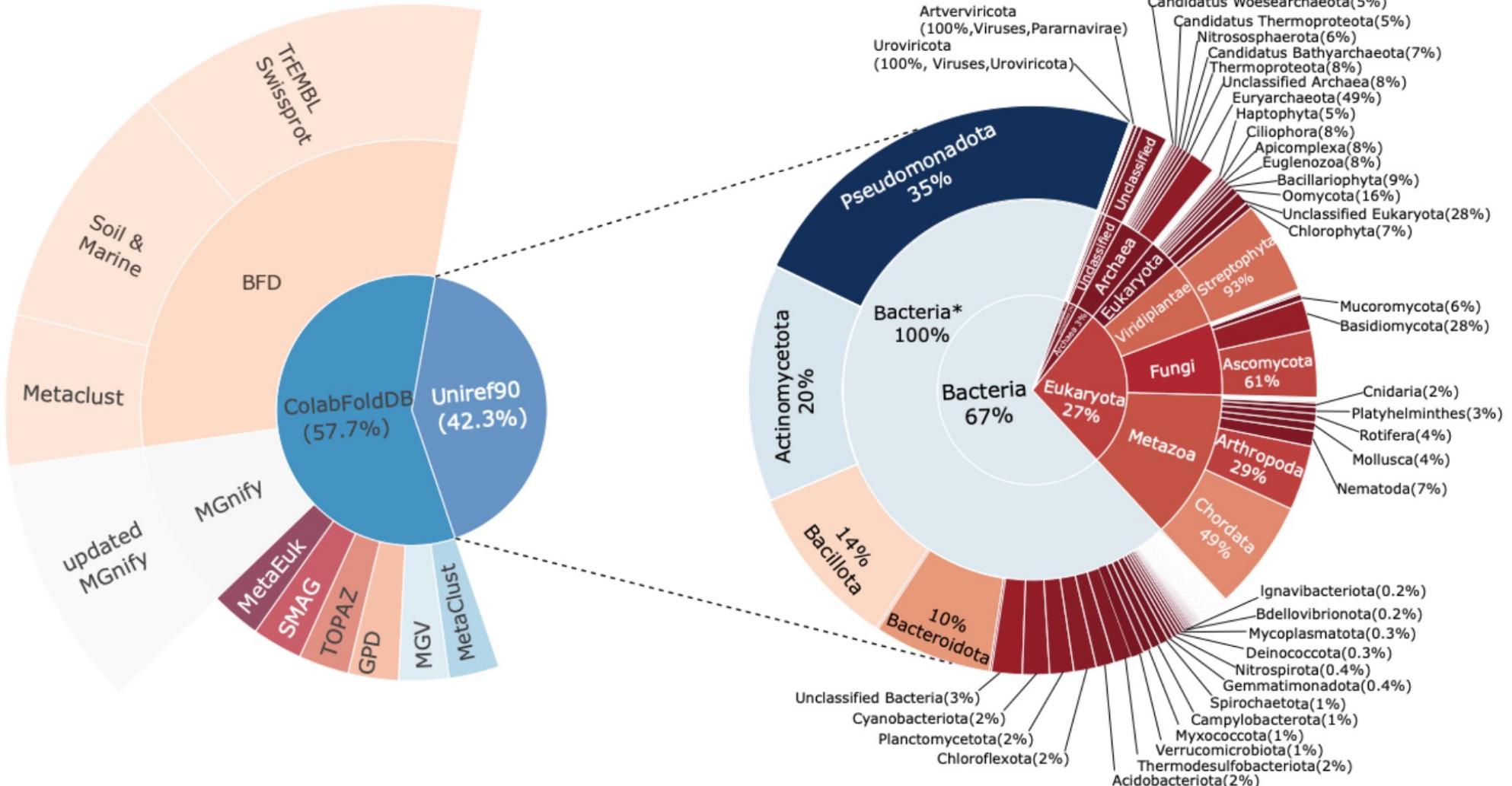
- Protein sequences are not i.i.d.
- Protein databases contain clusters of very similar sequences
- Models trained on all data overfit to largest cluster

# Deduplicate and cluster sequences for better performance



Depends on the task:

- Uniref50 may be better for **structure prediction**
- Greater number of sequences in Uniref90 is better for **mutation effect prediction**



# Most protein language models are trained on UniRef

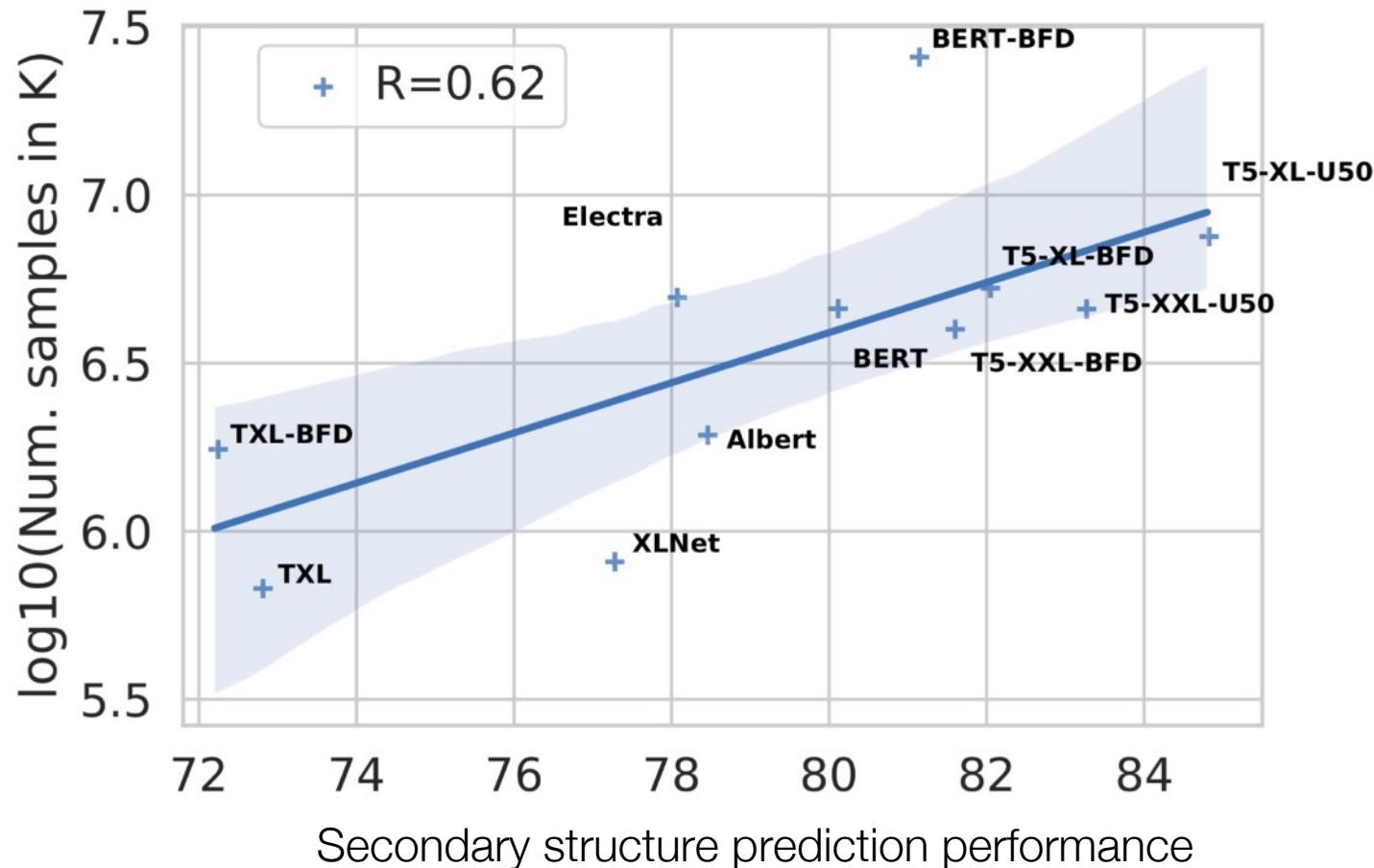
| Model and Repository  | Input | Network        | #Embedding | #Param. | Pretext Task  | Pre-training Dataset   |
|-----------------------|-------|----------------|------------|---------|---|------------------------|
| PRoBERTa [158]        | Seq   | Transformer    | 768        | 44M     | Masked Language Modeling                                | Swiss-Prot             |
| ESM-1b [14]           | Seq   | Transformer    | 1280       | 650M    | Masked Language Modeling                                | UniParc                |
| ProtTXL [12]          | Seq   | Transformer-XL | 1024       | 562M    | Masked Language Modeling                                | BFD100, UniRef100      |
| ProtBert [12]         | Seq   | BERT           | 1024       | 420M    | Masked Language Modeling                                | BFD100, UniRef100      |
| ProtXLNet [12]        | Seq   | XLNet          | 1024       | 409M    | Masked Language Modeling                                | UniRef100              |
| ProtAlbert [12]       | Seq   | ALBERT         | 4096       | 224M    | Masked Language Modeling                                | UniRef100              |
| ProtElectra [12]      | Seq   | ELECTRA        | 1024       | 420M    | Masked Language Modeling                                | UniRef100              |
| ProtT5 [12]           | Seq   | T5             | 1024       | 11B     | Masked Language Modeling                                | UniRef50, BFD100       |
| PMLM [159]            | Seq   | Transformer    | 1280       | 715M    | Masked Language Modeling                                | UniRef50               |
| MSA Transformer [160] | MSA   | Transformer    | 768        | 100M    | Masked Language Modeling                                | UniRef50, UniClust30   |
| ProteinLM [161]       | Seq   | BERT           | -          | 3B      | Masked Language Modeling                                | Pfam                   |
| PLUS-RNN [162]        | Seq   | RNN            | 2024       | 59M     | Masked Language Modeling<br>Same-Family Prediction      | Pfam                   |
| CARP [163]            | Seq   | CNN            | 1280       | 640M    | Masked Language Modeling                                | UniRef50               |
| AminoBERT [22]        | Seq   | Transformer    | 3072       | -       | Masked Language Modeling                                | UniParc                |
| OmegaPLM [164]        | Seq   | GAU [165]      | 1280       | 670M    | Masked Language Modeling<br>Span and Sequential Masking | UniRef50               |
| ProGen2 [166]         | Seq   | Transformer    | 4096       | 6.4B    | Masked Language Modeling                                | UniRef90, BFD30, BFD90 |
| ProtGPT2 [167]        | Seq   | GPT-2 [168]    | 1280       | 738M    | Next Amino Acid Prediction                              | UniRef50               |
| RITA [169]            | Seq   | GPT-3 [55]     | 2048       | 1.2B    | Next Amino Acid Prediction                              | UniRef100              |
| ESM-2 [13]            | Seq   | Transformer    | 5120       | 15B     | Masked Language Modeling                                | UniRef50               |

# What are the sizes of these datasets?

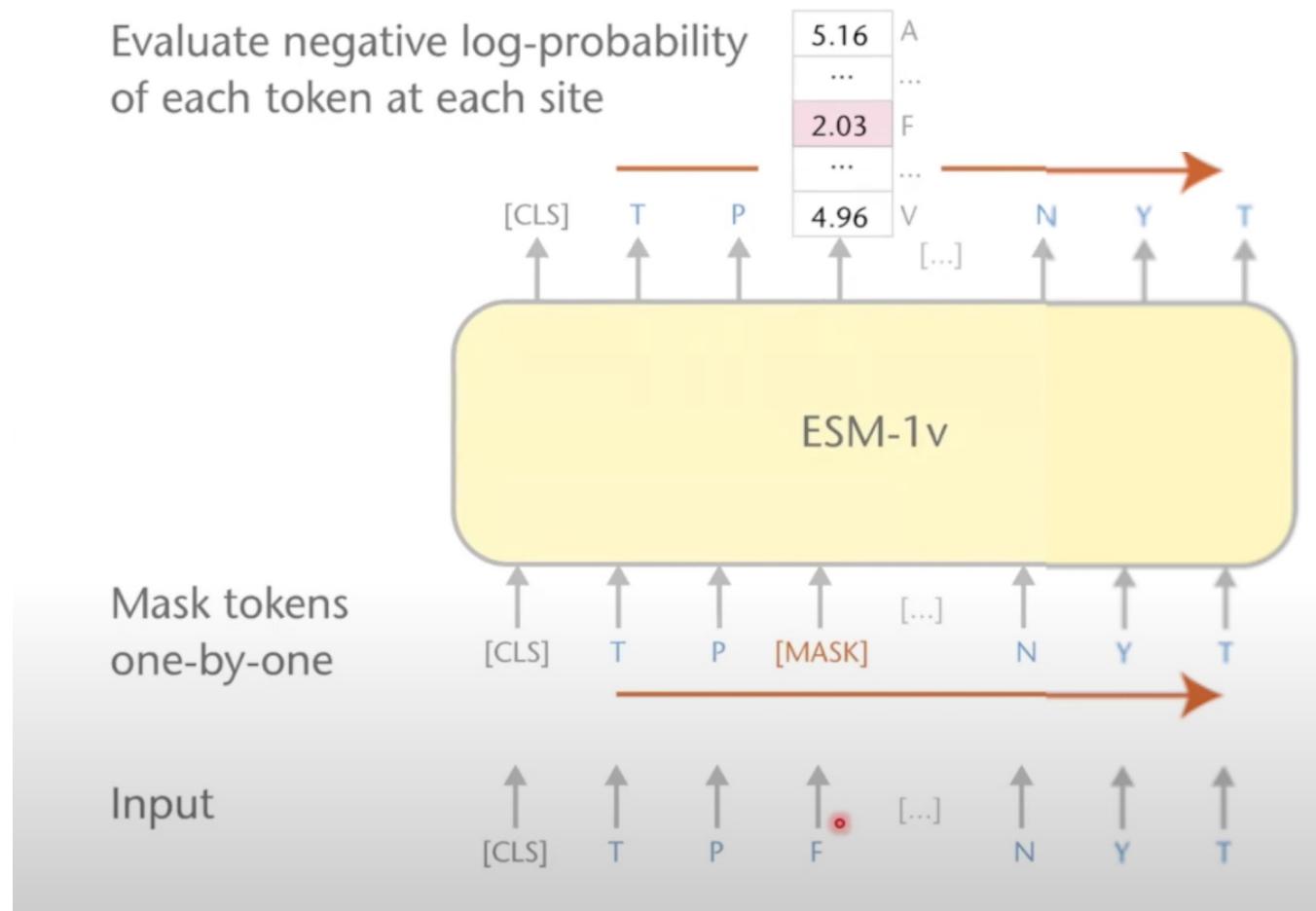
| Dataset                    | #Proteins | Disk Space | Description                                      |
|----------------------------|-----------|------------|--|
| UniProtKB/Swiss-Prot [9]   | 500K      | 0.59GB     | knowledgebase                                    |
| UniProtKB/TrEMBL [9]       | 229M      | 146GB      | knowledgebase                                    |
| UniRef100 [215]            | 314M      | 76.9GB     | clustered sets of sequences                      |
| UniRef90 [215]             | 150M      | 34GB       | 90% identity                                     |
| UniRef50 [215]             | 53M       | 10.3GB     | 50% identity                                     |
| UniParc [9]                | 528M      | 106GB      | sequence   |
| PDB [8]                    | 180K      | 50GB       | 3D structure                                     |
| CATH4.3 [216]              | -         | 1073MB     | hierarchical classification                      |
| BFD [217]                  | 2500M     | 272GB      | sequence profile                                 |
| Pfam [218]                 | 47M       | 14.1GB     | protein families                                 |
| AlphaFoldDB [219]          | 214M      | 23 TB      | predicted 3D structures                          |
| ESM Metagenomic Atlas [13] | 772M      | -          | predicted metagenomic protein structures         |
| ColAbFoldDB [170]          | 950M      | -          | an amalgamation of various metagenomic databases |
| ProteinKG25 [10]           | 5.6M      | 147MB      | a knowledge graph dataset with GO                |
| Uniclust30 [220]           | -         | 6.6GB      | clustered protein sequences                      |
| SCOP [221]                 | -         | -          | structural classification                        |
| SCOPe [222]                | -         | 86MB       | an extended version of SCOP                      |
| OpenProteinSet [223]       | 16M       | -          | MSAs   |

BFD is the biggest (and most fantastic)

# Model performance improves with more training sequences



# Mutation effects from a protein language model



Compute **Log Ratio** for each mutant

$$\log \frac{p(\mathbf{x}_{\text{mut}} | \theta)}{p(\mathbf{x}_{\text{wild}} | \theta)}$$

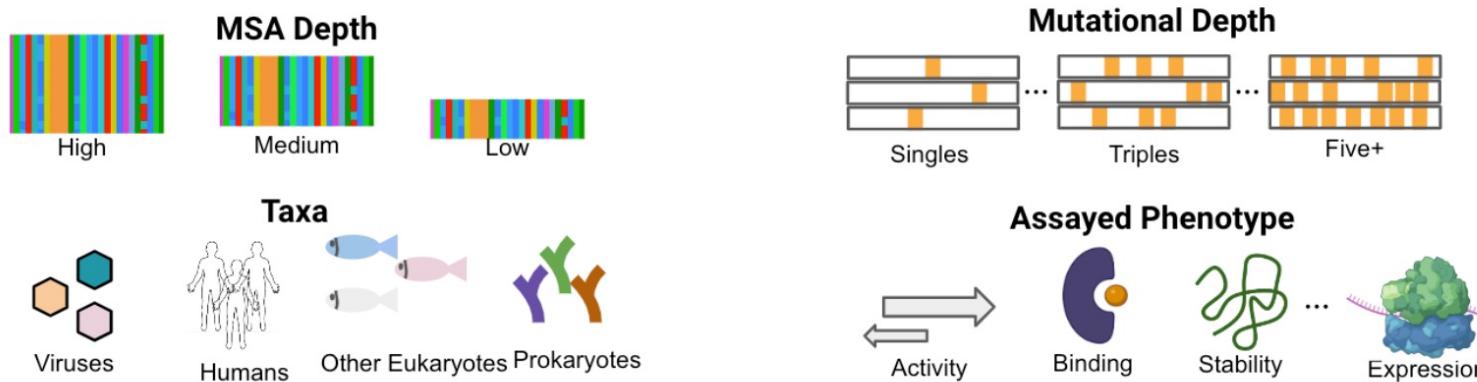
# Common evaluation tasks

- **Contact Prediction:** predict whether two residues are in the pre-defined proximity
- **Fold Prediction:** assign sequence to one to 1,195 known folds
- **Secondary Structure Prediction:** classify sequence into Helix, Strand, Coil
- **Solubility:** binary
- **Stability:** whether protein remains folded
- **Localization:** distribution into 100 Unique Subcellular categories
- **TCR-pMHC Affinity:** interaction between T cell receptors and peptide-major histocompatibility complex
- ....

# How to evaluate mutation effect prediction across all of the models?

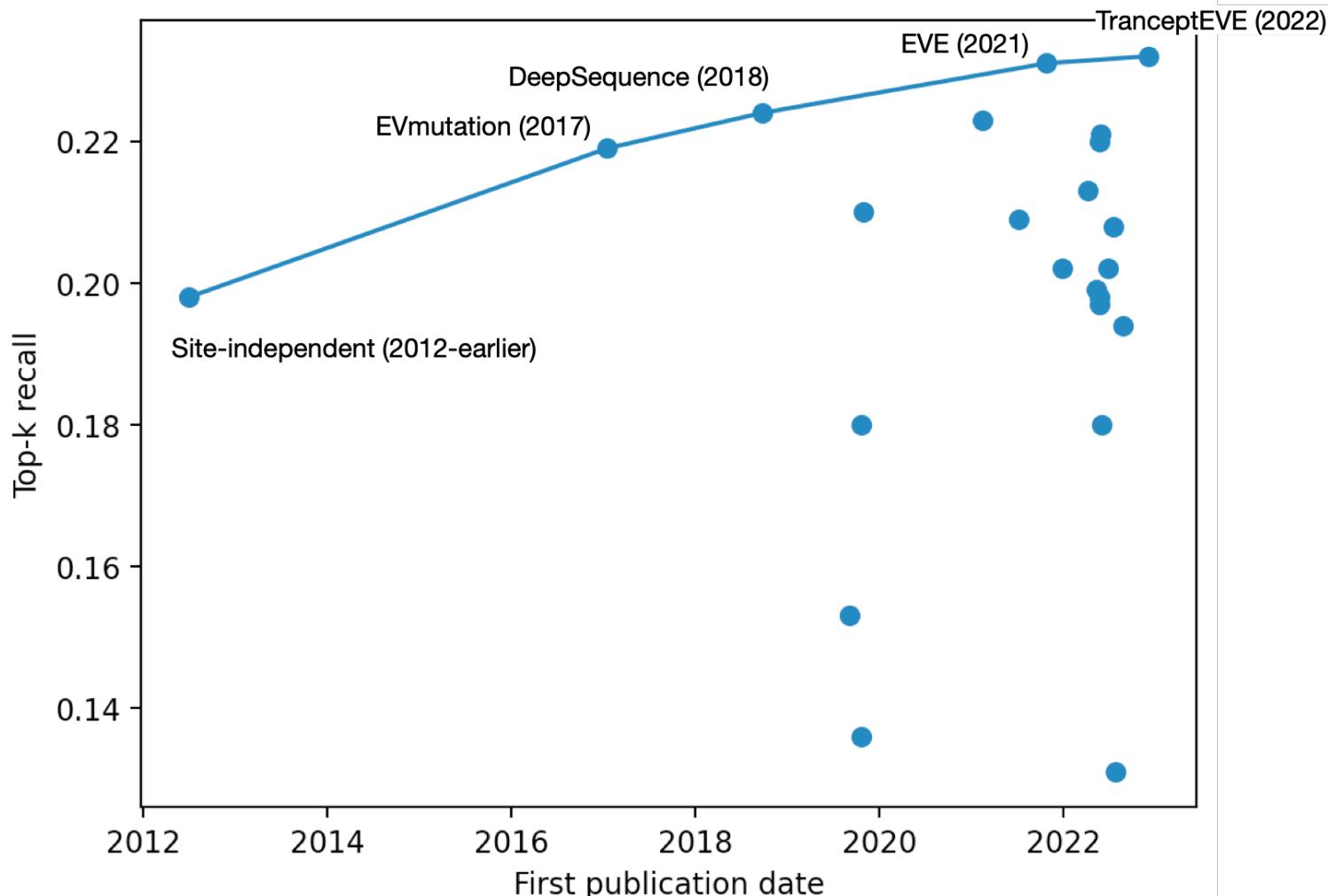
**ProteinGym covers a diverse array of proteins**

~30 proteins (2018) → >200 proteins (now)

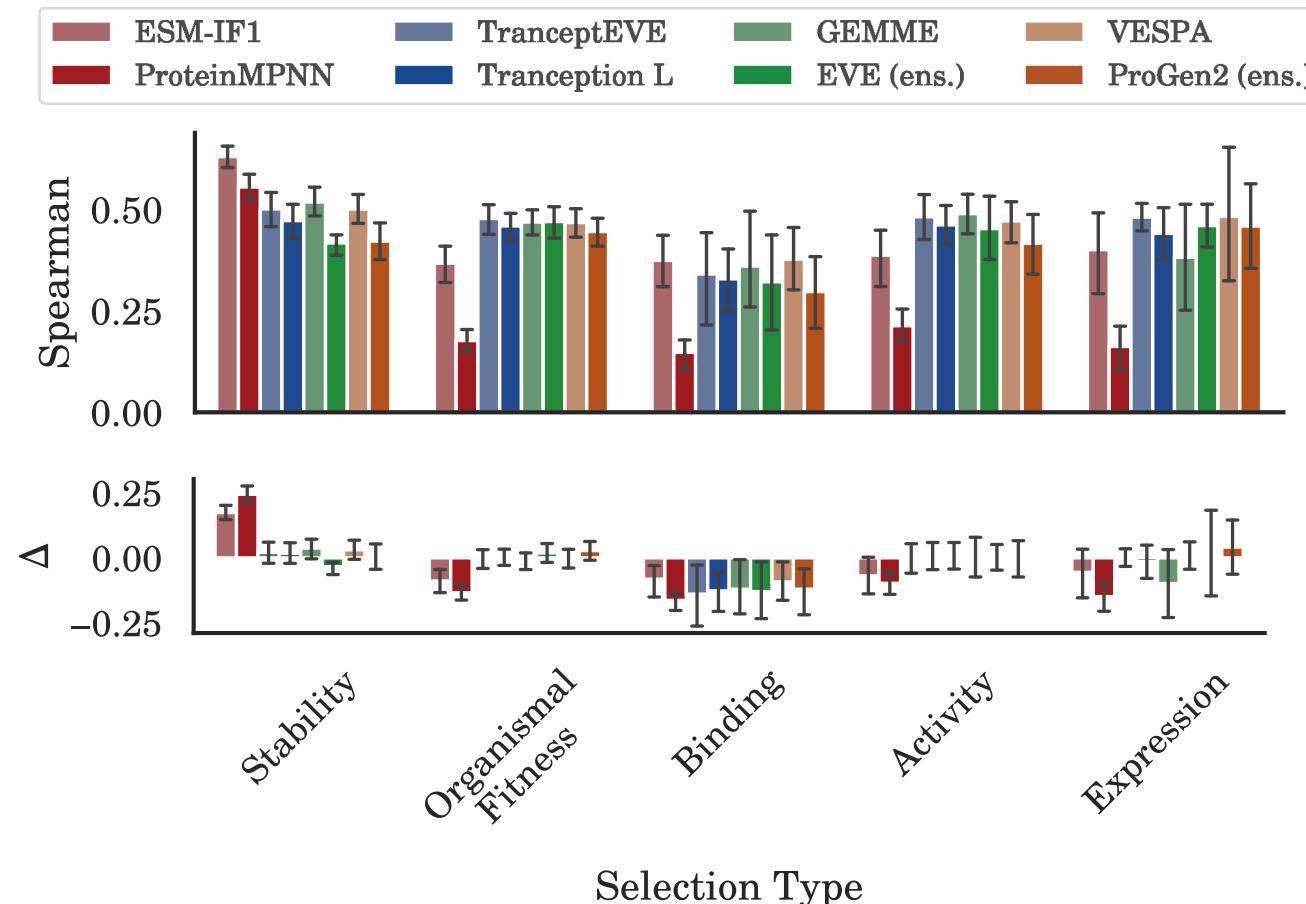


Total: 2.8 million variants

# Performance against Deep Mutational Scans has improved over time



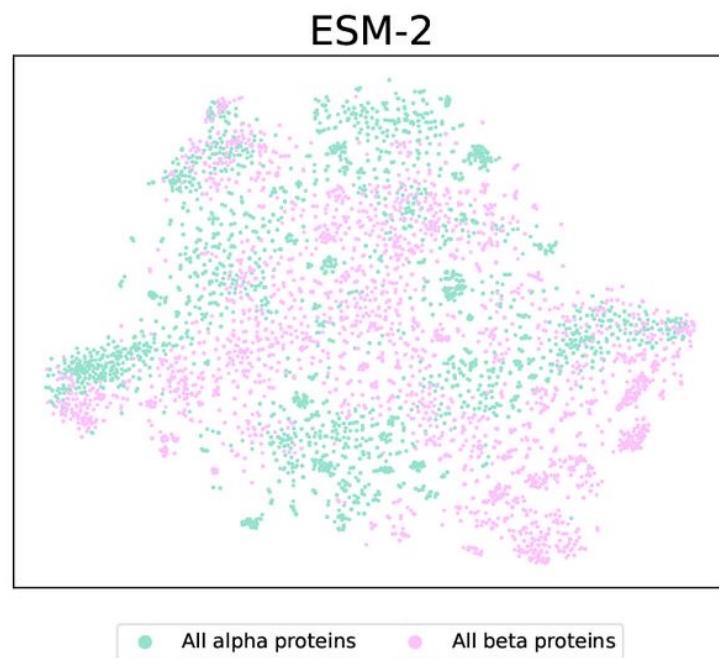
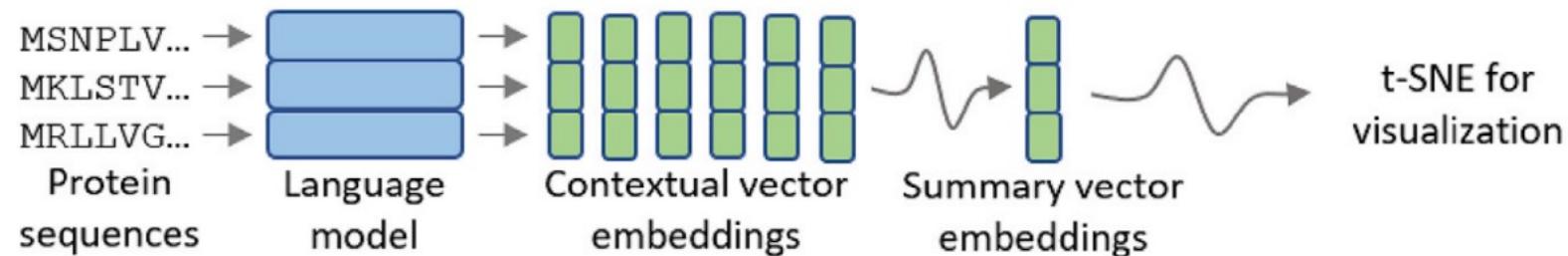
# Different models better at predicting different functions



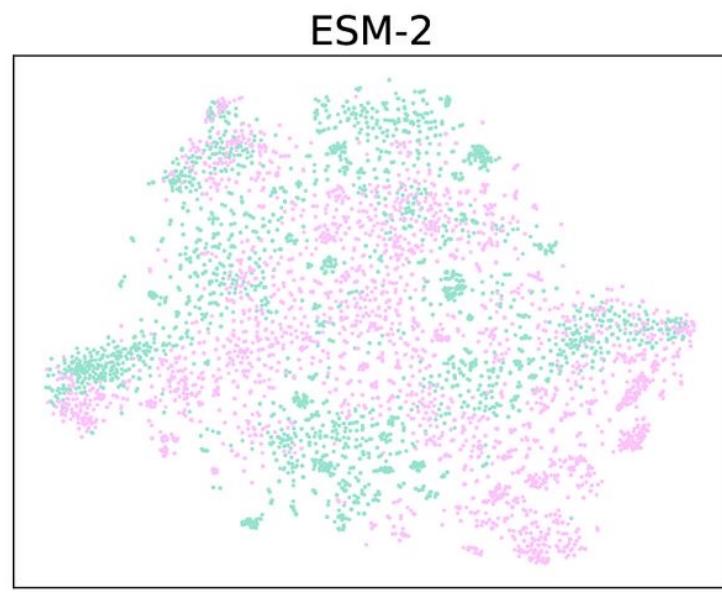
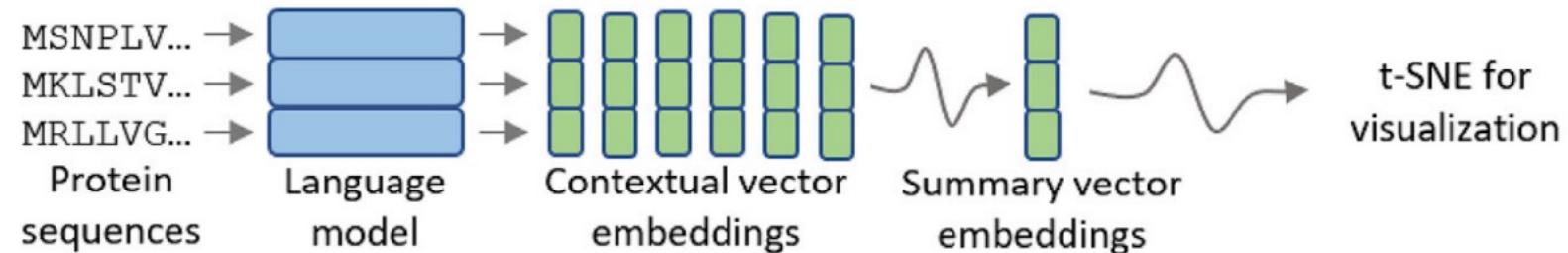
# Models do better with more relevant sequences

| Model type       | Model name             | Spearman by MSA depth (↑) |              |              |              |
|------------------|------------------------|---------------------------|--------------|--------------|--------------|
|                  |                        | Low                       | Medium       | High         | All          |
| Alignment-based  | Site-Independent       | 0.426                     | 0.373        | 0.320        | 0.373        |
|                  | WaveNet                | 0.207                     | 0.255        | 0.207        | 0.223        |
|                  | EVmutation             | 0.403                     | 0.423        | 0.410        | 0.412        |
|                  | DeepSequence (ens.)    | 0.383                     | 0.428        | 0.473        | 0.428        |
|                  | EVE (ens.)             | 0.425                     | 0.453        | 0.481        | 0.453        |
|                  | GEMME                  | <b>0.455</b>              | <b>0.470</b> | 0.497        | <b>0.474</b> |
| Protein language | UniRep                 | 0.181                     | 0.161        | 0.209        | 0.184        |
|                  | CARP (640M)            | 0.314                     | 0.375        | 0.428        | 0.372        |
|                  | ESM-1b                 | 0.350                     | 0.398        | 0.482        | 0.410        |
|                  | ESM-2 (15B)            | 0.357                     | 0.414        | 0.473        | 0.415        |
|                  | RITA XL                | 0.315                     | 0.382        | 0.412        | 0.370        |
|                  | ESM-1v (ens.)          | 0.326                     | 0.418        | 0.502        | 0.415        |
|                  | ProGen2 XL             | 0.354                     | 0.405        | 0.444        | 0.401        |
|                  | VESPA                  | 0.427                     | 0.455        | 0.484        | 0.455        |
| Hybrid           | UniRep evotuned        | 0.330                     | 0.344        | 0.372        | 0.349        |
|                  | MSA Transformer (ens.) | 0.404                     | 0.450        | 0.488        | 0.447        |
|                  | Tranception L          | 0.432                     | 0.438        | 0.473        | 0.448        |
|                  | TranceptEVE L          | 0.451                     | 0.467        | 0.492        | 0.470        |
| Inverse Folding  | ESM-IF1                | 0.300                     | 0.431        | <b>0.544</b> | 0.425        |
|                  | MIF-ST                 | 0.376                     | 0.403        | 0.456        | 0.412        |
|                  | ProteinMPNN            | 0.173                     | 0.280        | 0.434        | 0.296        |

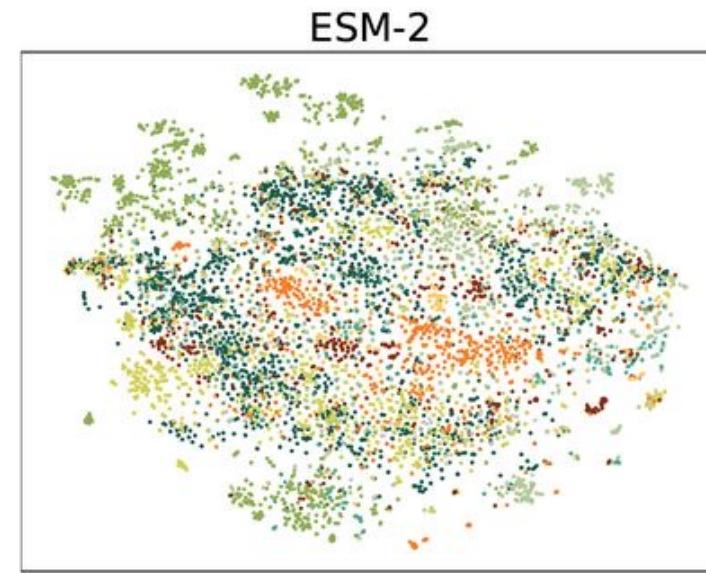
# ESM embedding space



# ESM embedding space



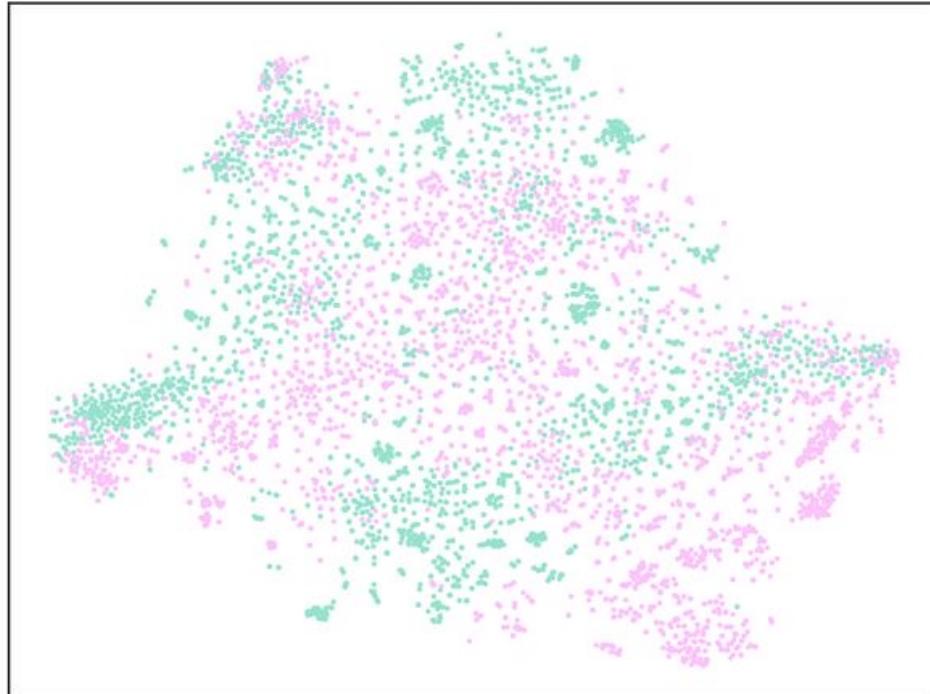
All alpha proteins      All beta proteins



|                 |                         |
|-----------------|-------------------------|
| ● Nucleus       | ● Endoplasmic reticulum |
| ● Cytoplasm     | ● Golgi apparatus       |
| ● Extracellular | ● Chloroplast           |
| ● Mitochondrion | ● Lysosome              |
| ● Cell.membrane | ● Peroxisome            |

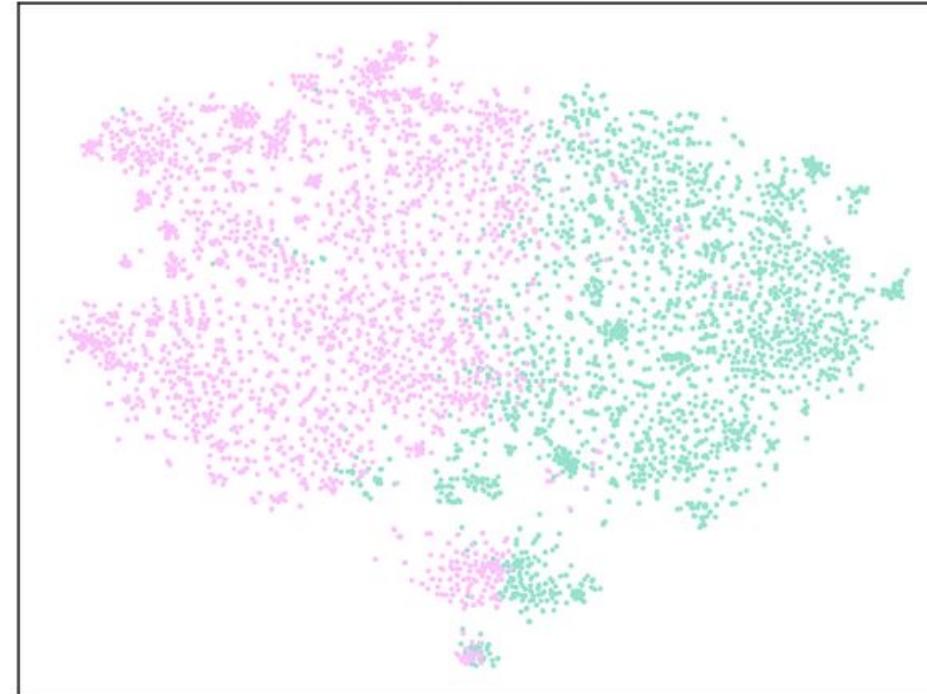
# Structure-aware tokens separate embeddings by structure type

ESM-2



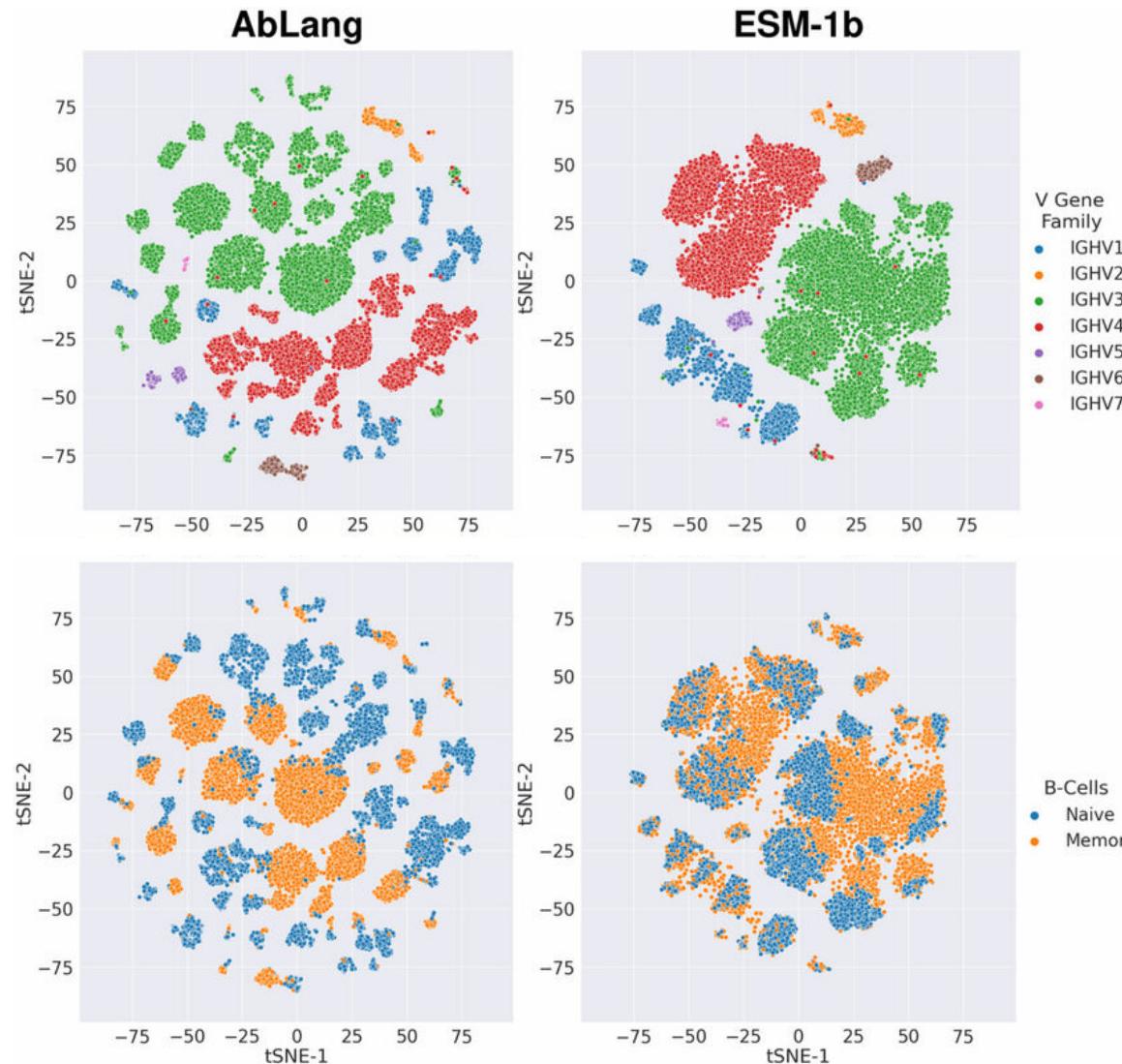
● All alpha proteins ● All beta proteins

SaProt

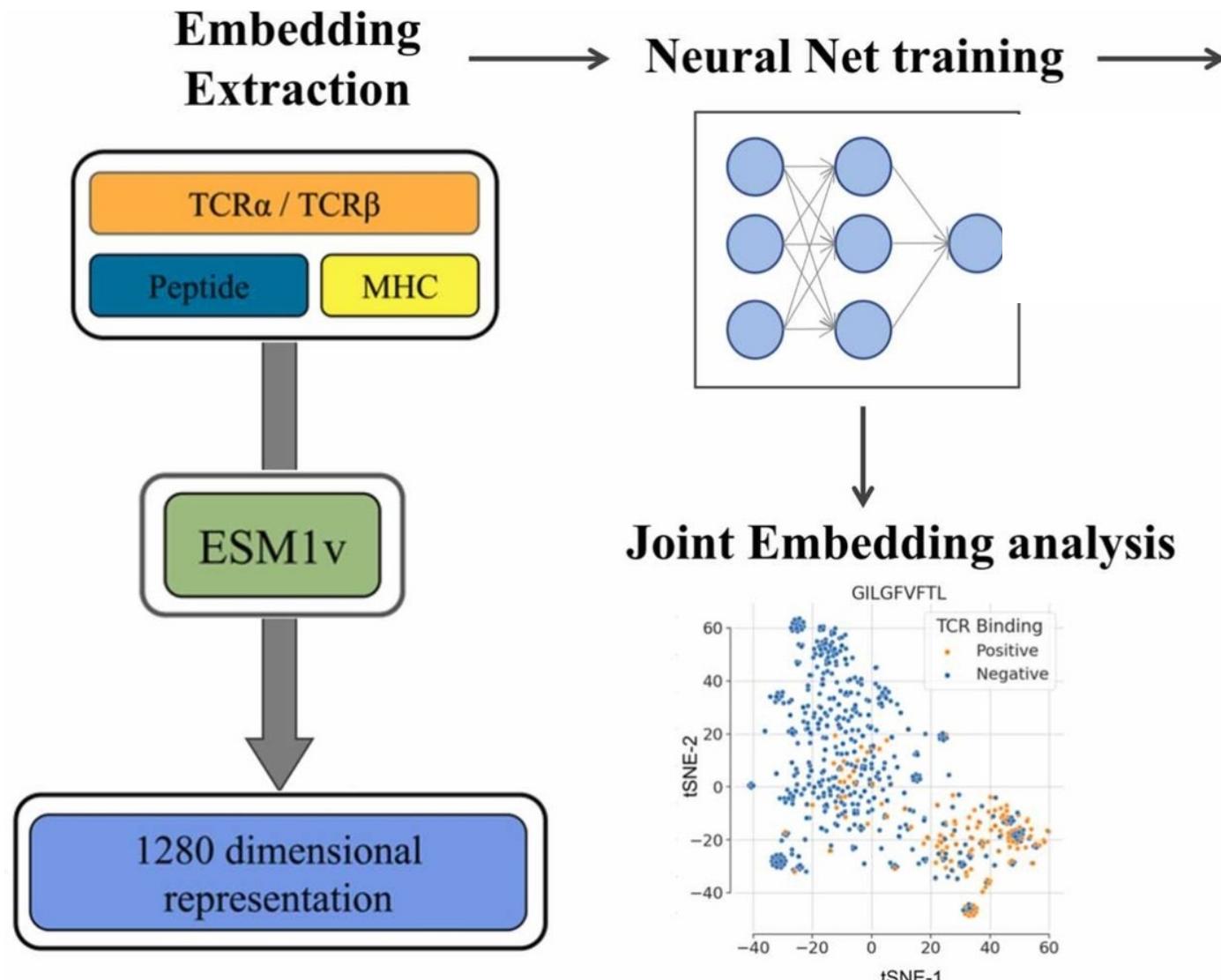


● All alpha proteins ● All beta proteins

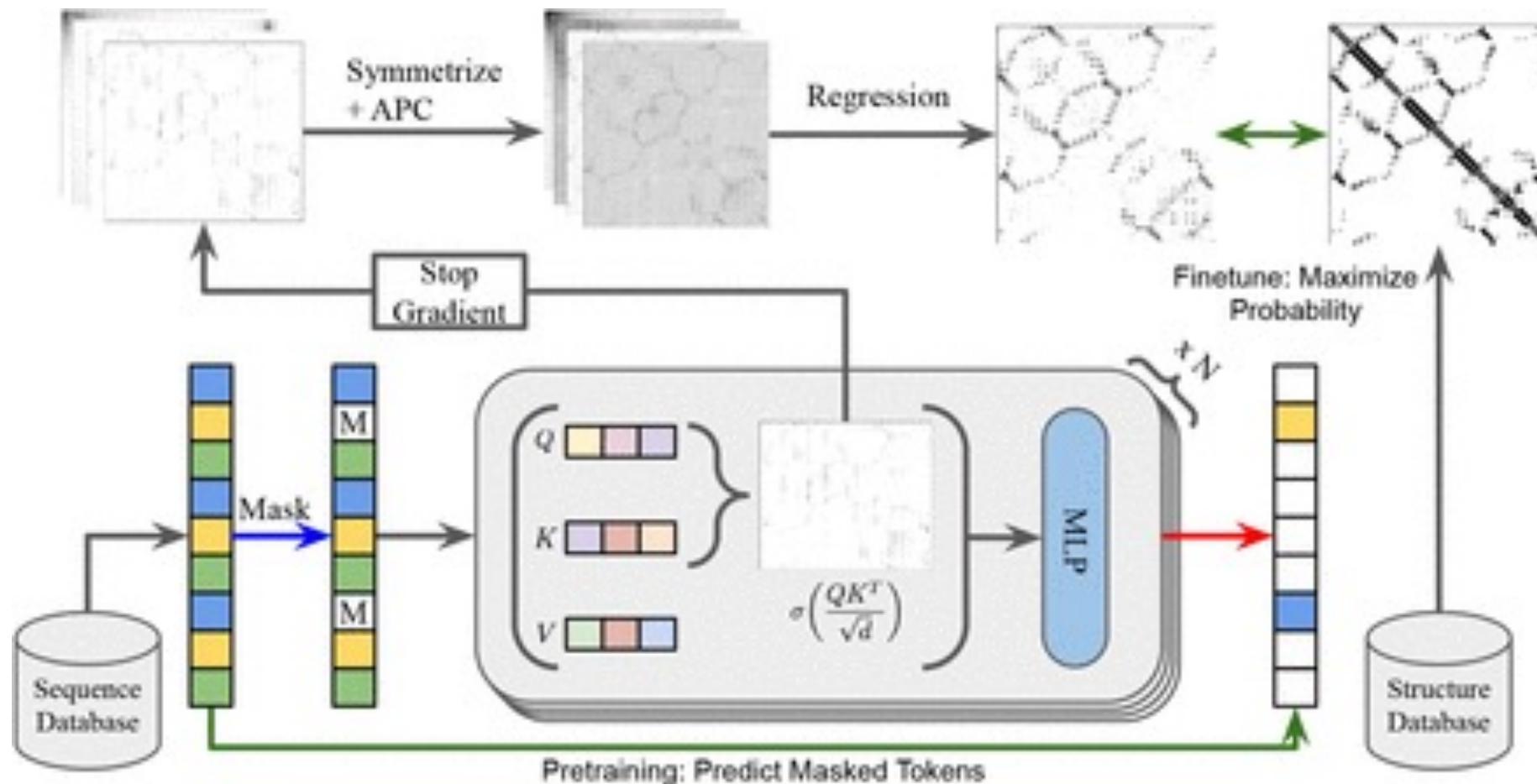
# Focusing on one class of proteins (i.e., antibodies) leads to informative embeddings for that class

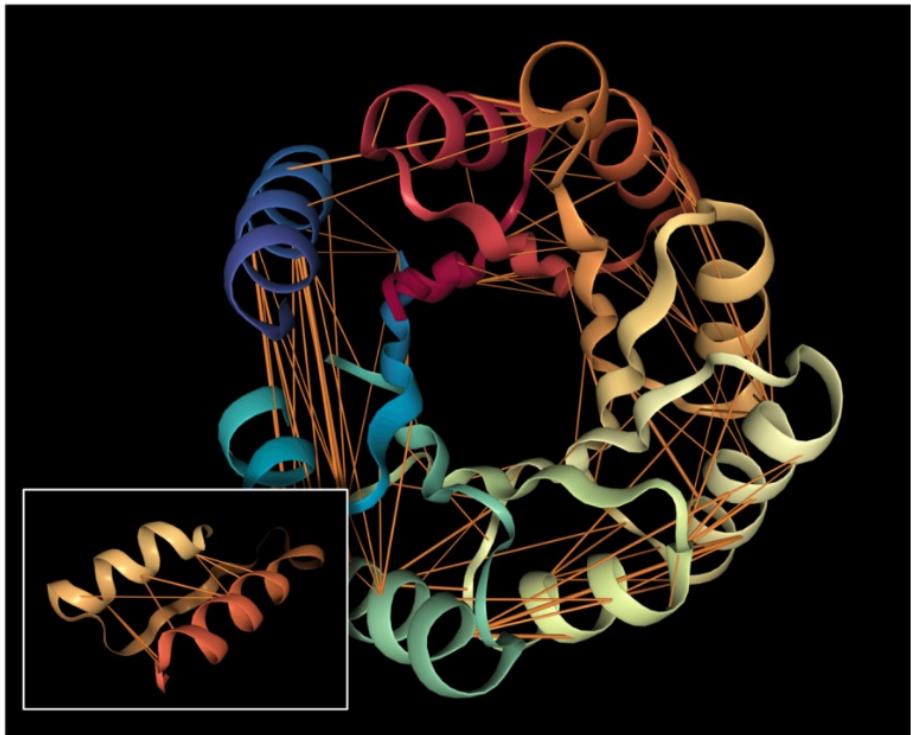


# Embeddings as inputs into supervised models

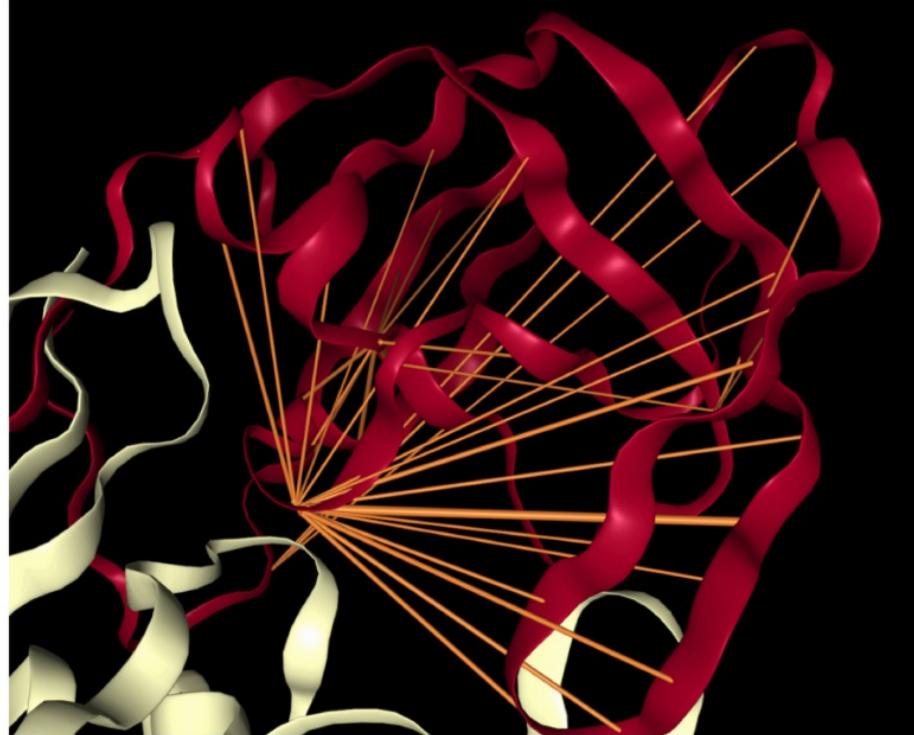


# Language model attention maps learn structural contacts



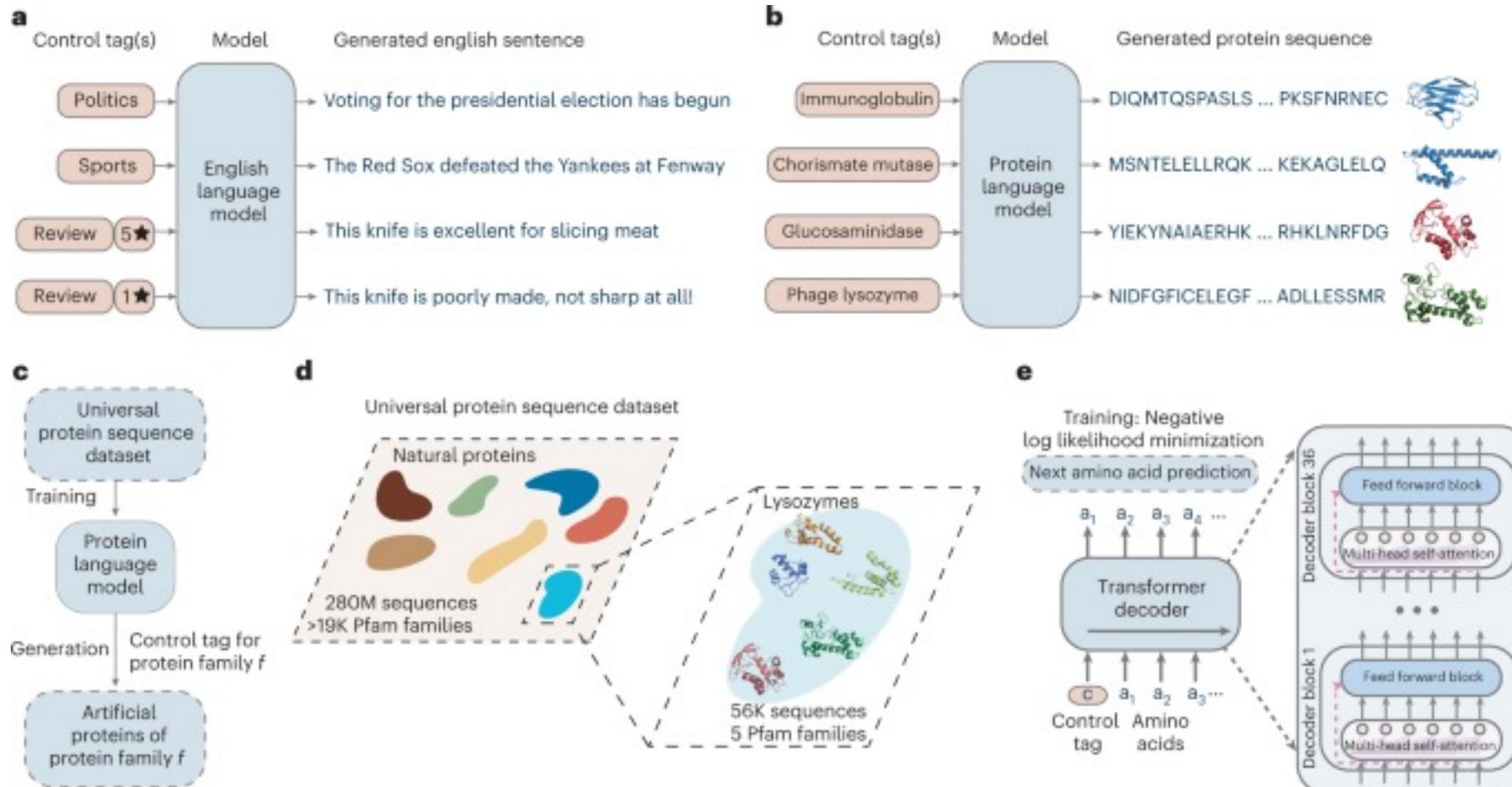


(a) Attention in head 12-4, which targets amino acid pairs that are close in physical space (see inset subsequence 117D-157I) but lie apart in the sequence. Example is a *de novo* designed TIM-barrel (5BVL) with characteristic symmetry.

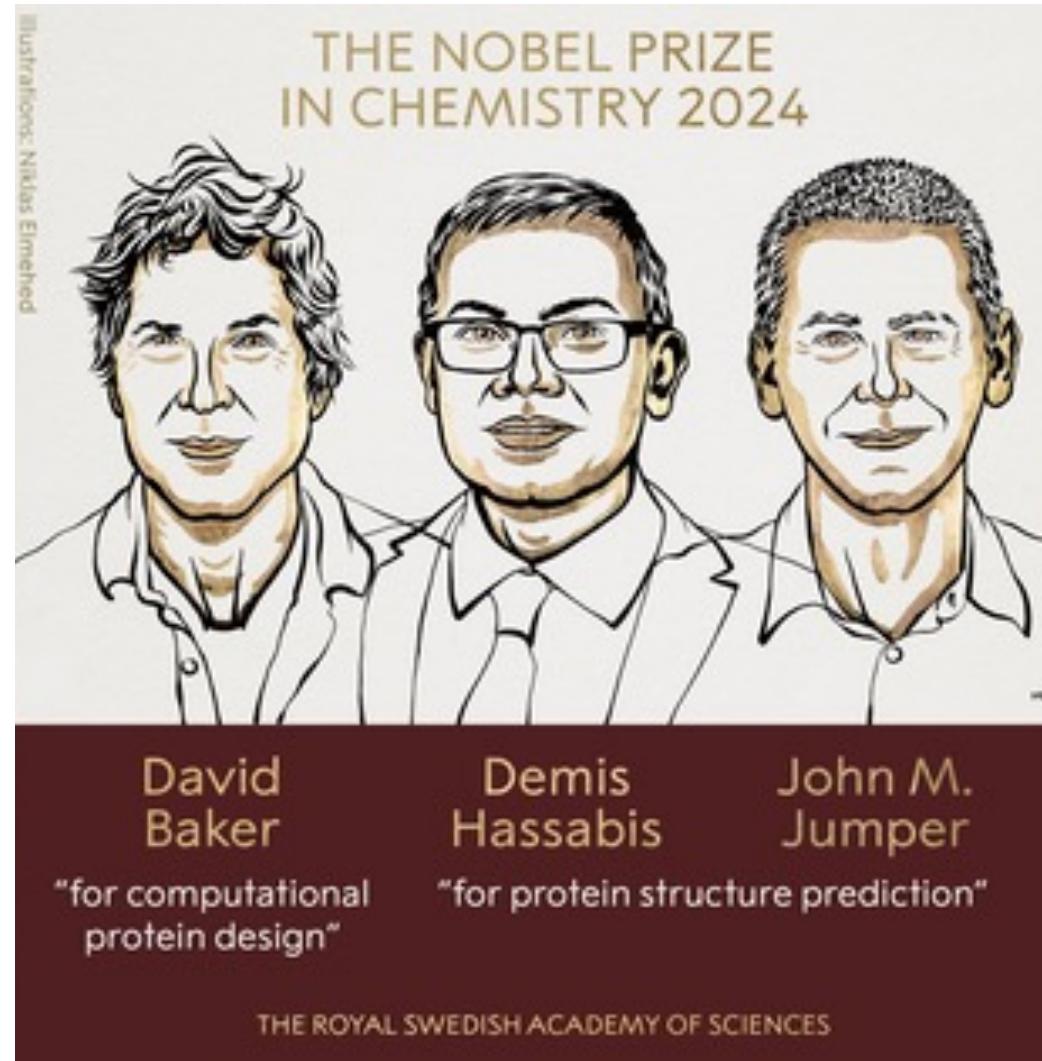


(b) Attention in head 7-1, which targets binding sites, a key functional component of proteins. Example is HIV-1 protease (7HVP). The primary location receiving attention is 27G, a binding site for protease inhibitor small-molecule drugs.

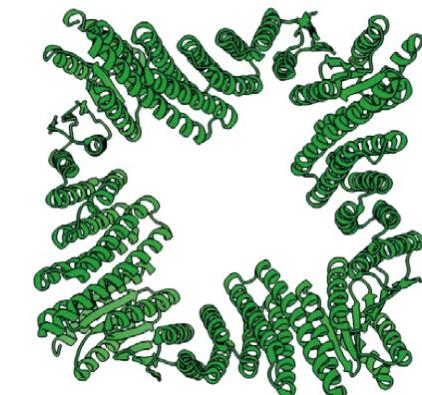
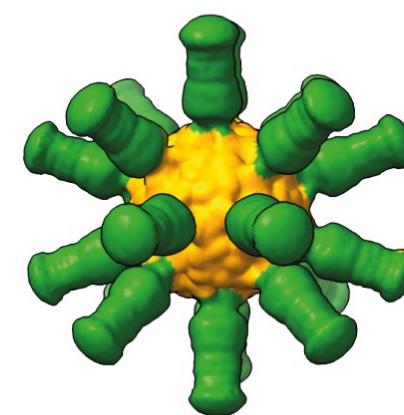
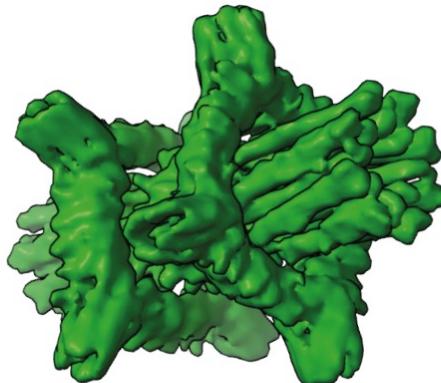
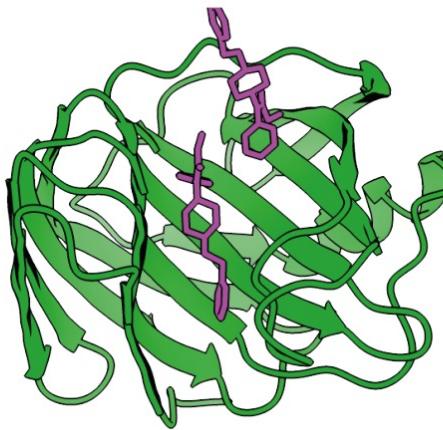
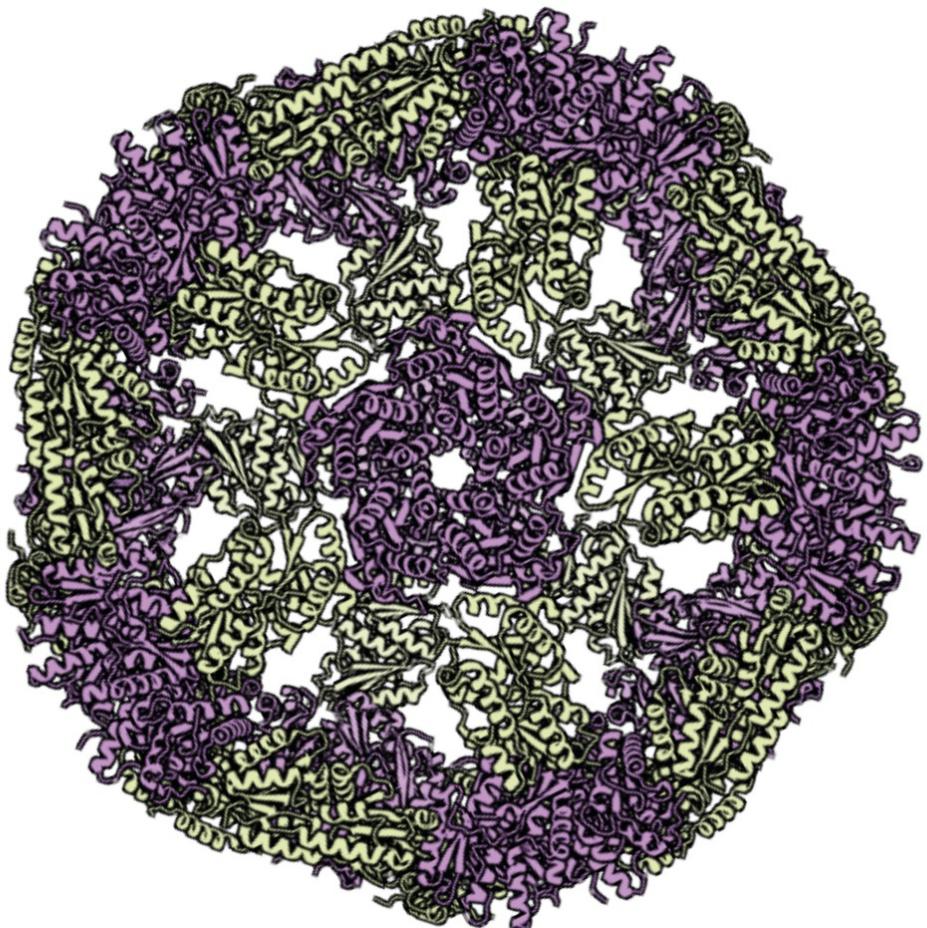
# Protein language models can design new proteins!



# An exciting time for computational protein design!



**Figure 4.** Proteins developed using Baker's program Rosetta.



# A case study on predicting viral evolution of the SARS-CoV-2 Spike protein

Pandemic onset

Dec 2019/ Jan 2020



Feb. 2020

Aug.

Feb. 2021

Aug.

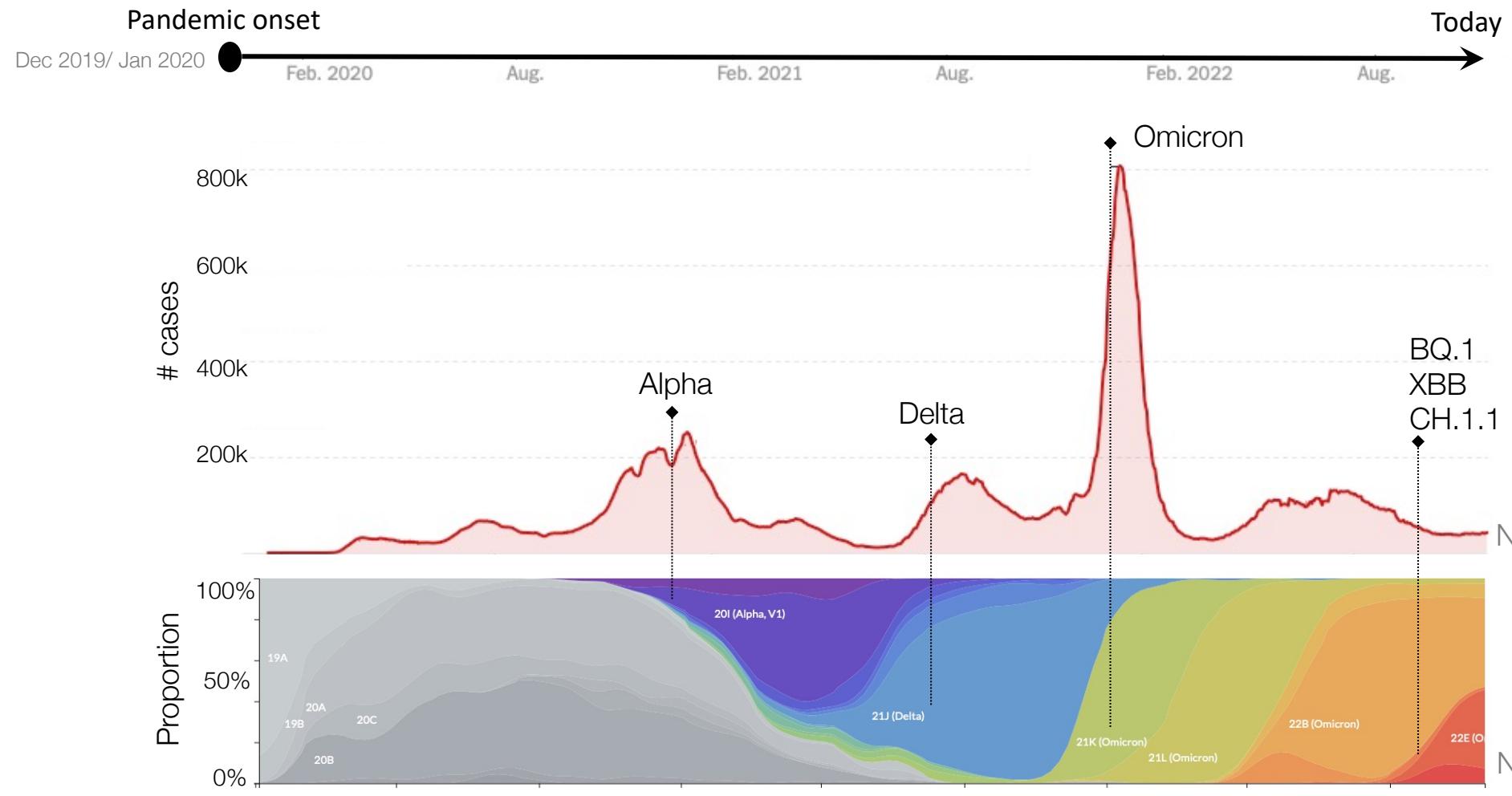
Feb. 2022

Aug.

Today



NY Times 11/25/2022



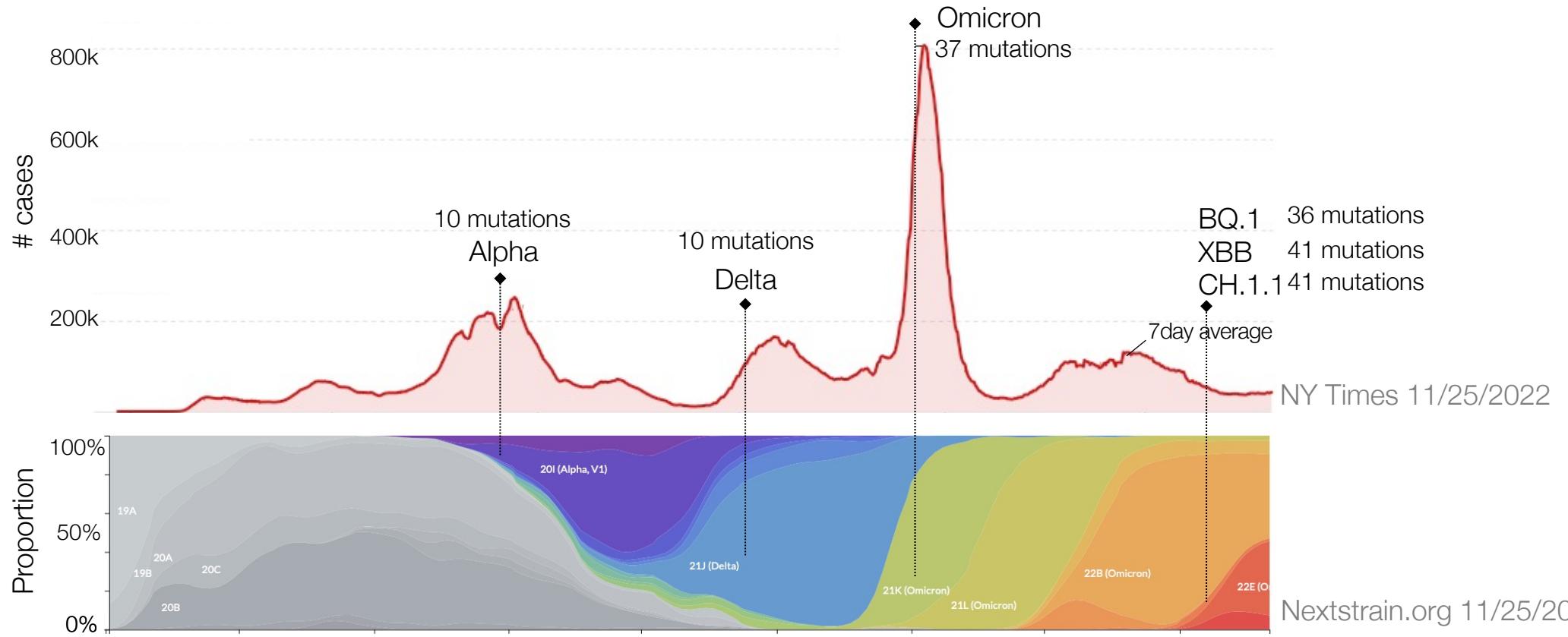
**HEALTH**  
**The coronavirus isn't mutating quickly, suggesting a vaccine would offer lasting protection**



By Joel Achenbach

March 24, 2020 at 4:30 p.m. EDT

Washington post – March 2020



HEALTH

The coronavirus isn't mutating quickly, suggesting a vaccine would offer lasting protection

By Joel Achenbach  
March 24, 2020 at 4:30 p.m. EDT

Washington post – March 2020

Caused significant antibody escape

Can machine learning facilitate ***early*** and  
***accurate*** predictions of viral escape from  
antibodies?

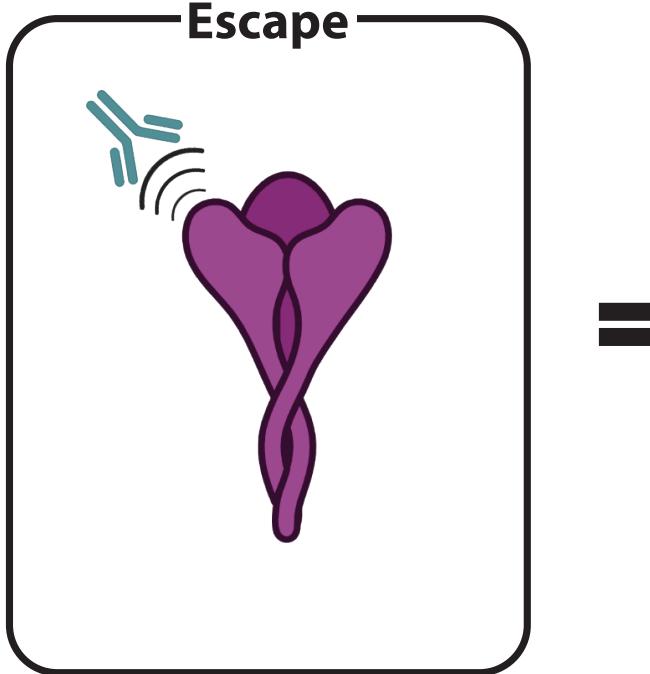
# A trip back in time to the available data for modeling



>15 million SARS-CoV-2 genomes  
>690 Spike-Ab Structures  
many experiments

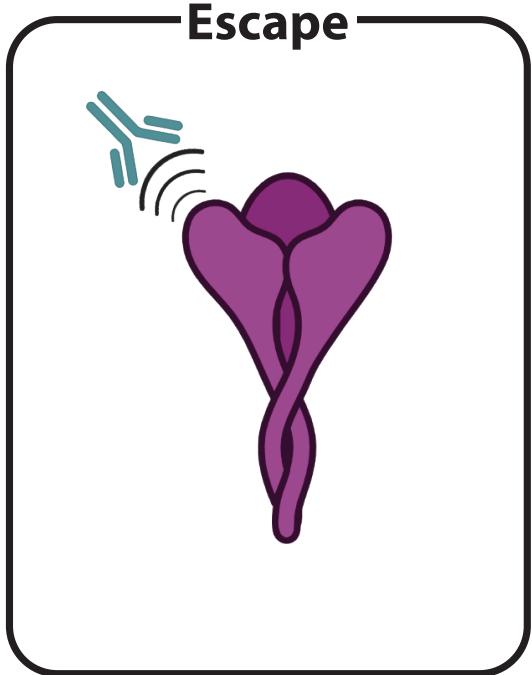
Not using any data after  
Jan 2020 for model  
training

# How?

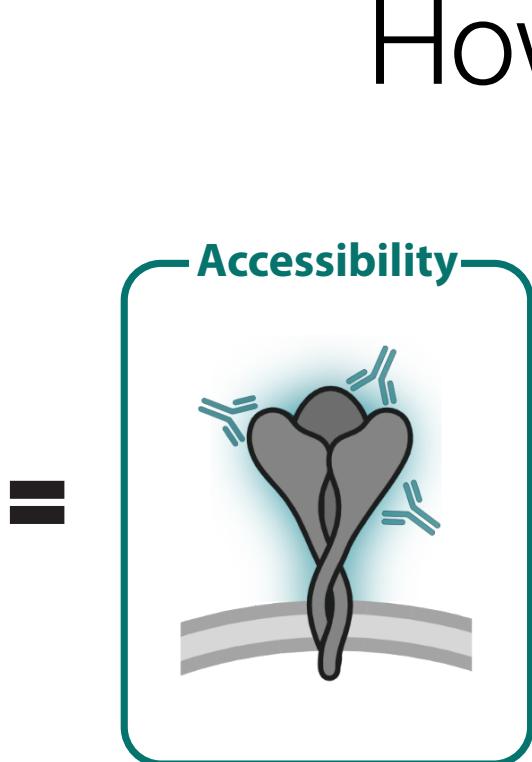


Probability that  
mutant escapes  
antibody

# How?

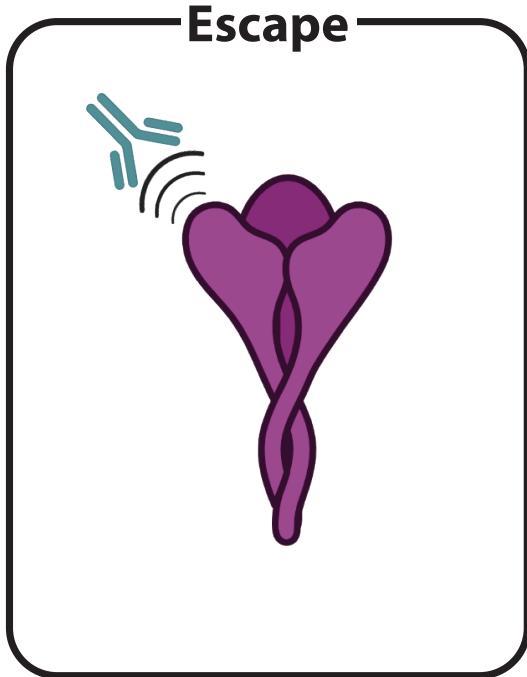


Probability that  
mutant escapes  
antibody

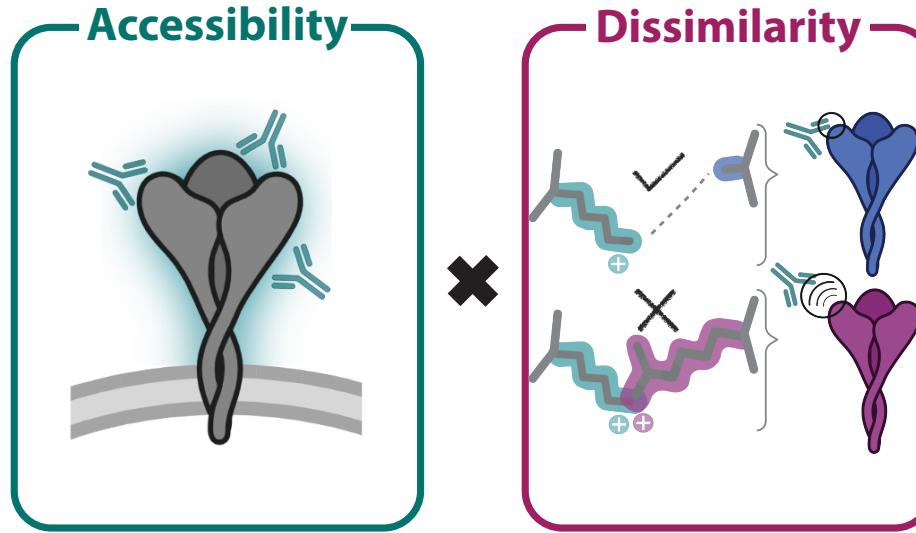


Mutated residue is  
accessible to antibody

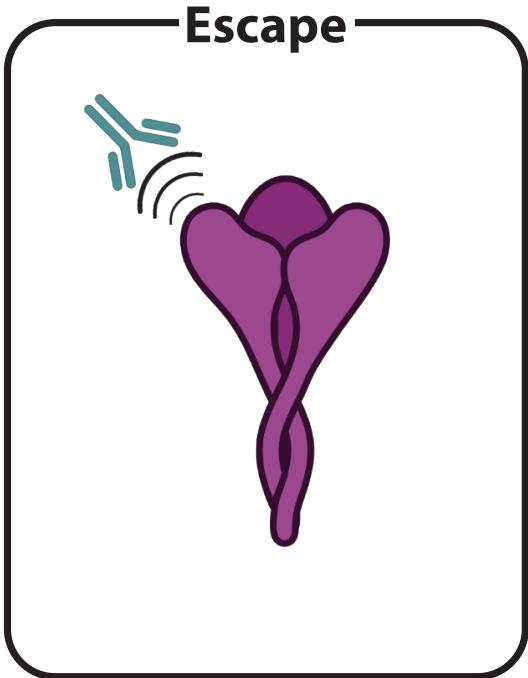
# How?



=

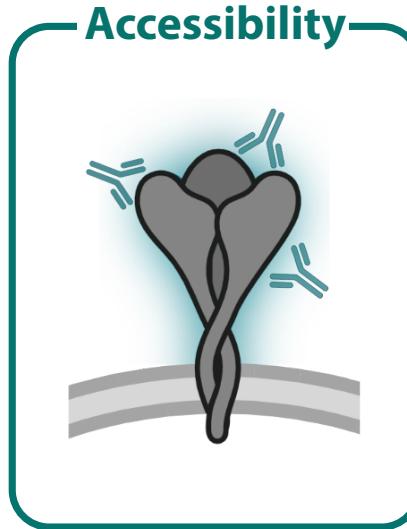


# How?



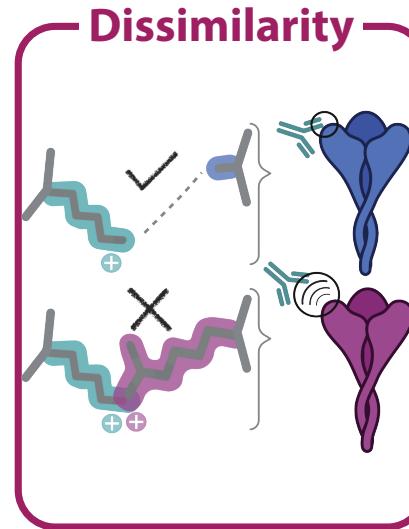
Probability that mutant  
escapes antibody

=



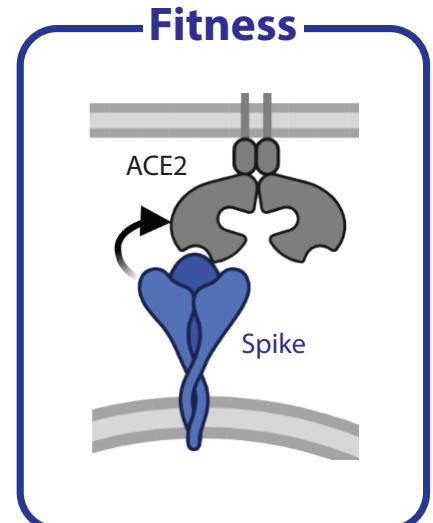
Mutated residue is  
accessible to antibody

\*



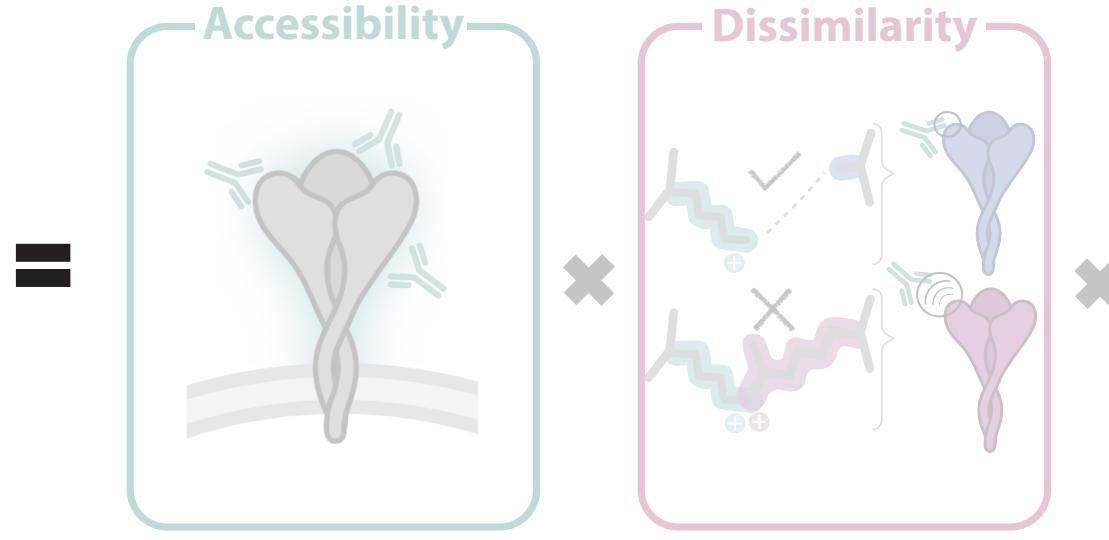
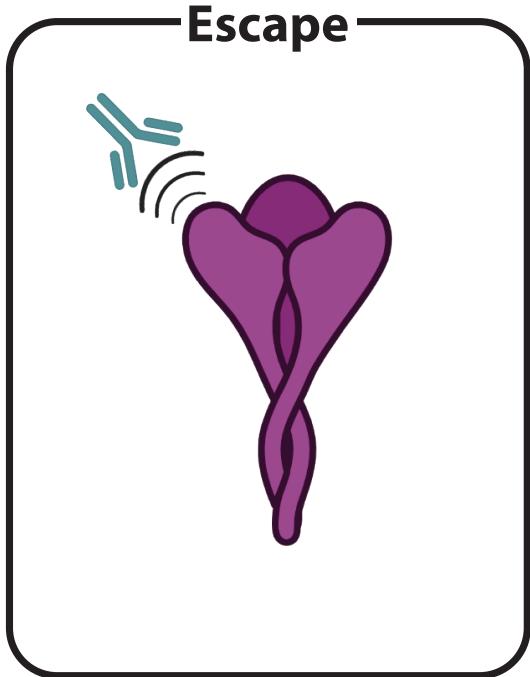
Mutation disrupts  
antibody binding

\*

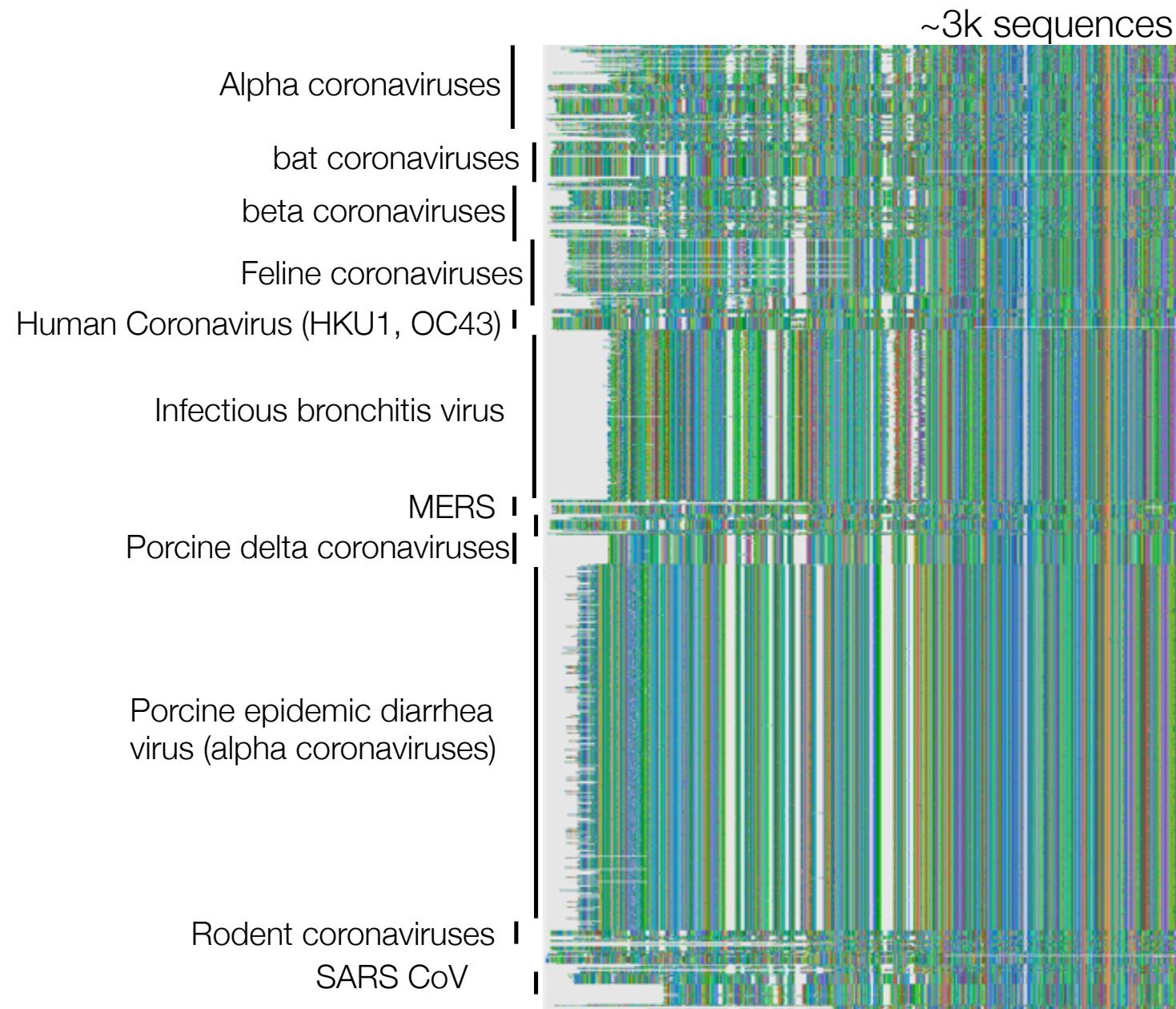
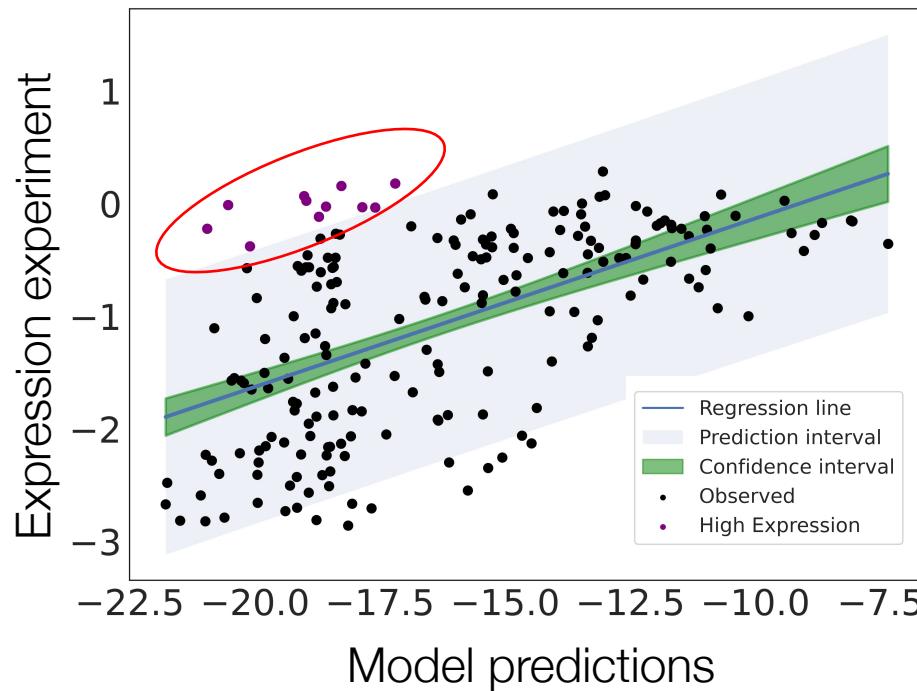


Mutation maintains  
viral fitness

# How?

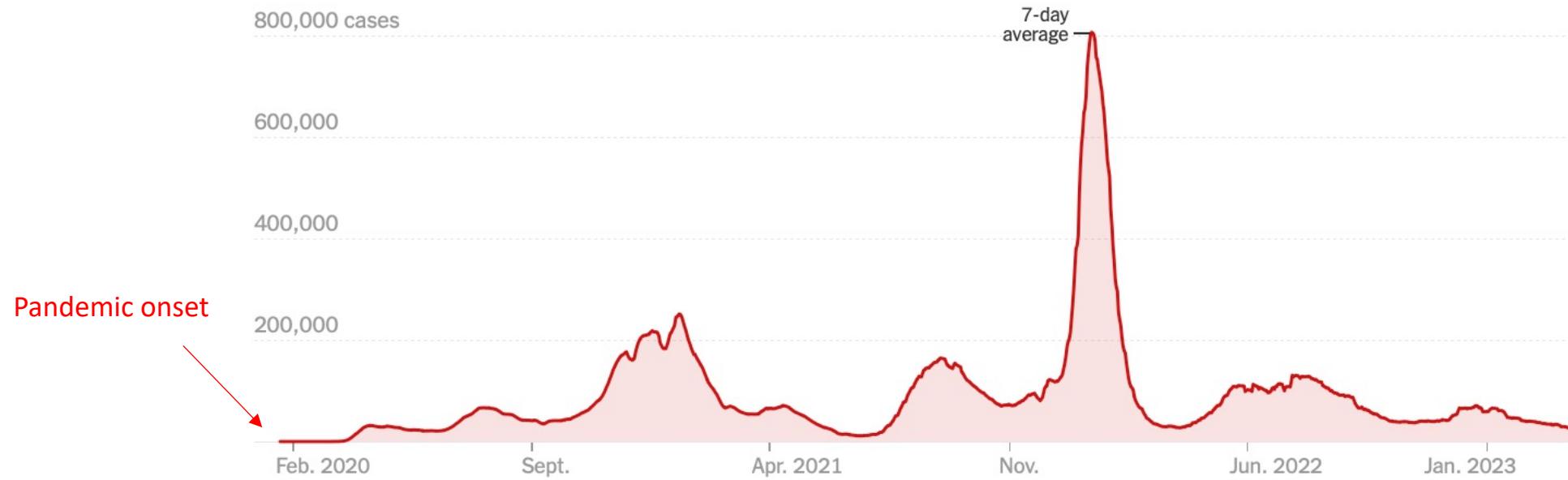


# Deep learning model trained on pre-2020 data predicts SARS-CoV-2 fitness



Can EVEscape trained  
on pre-pandemic data *anticipate pandemic  
mutations?*

# A trip back in time to evaluate predictions



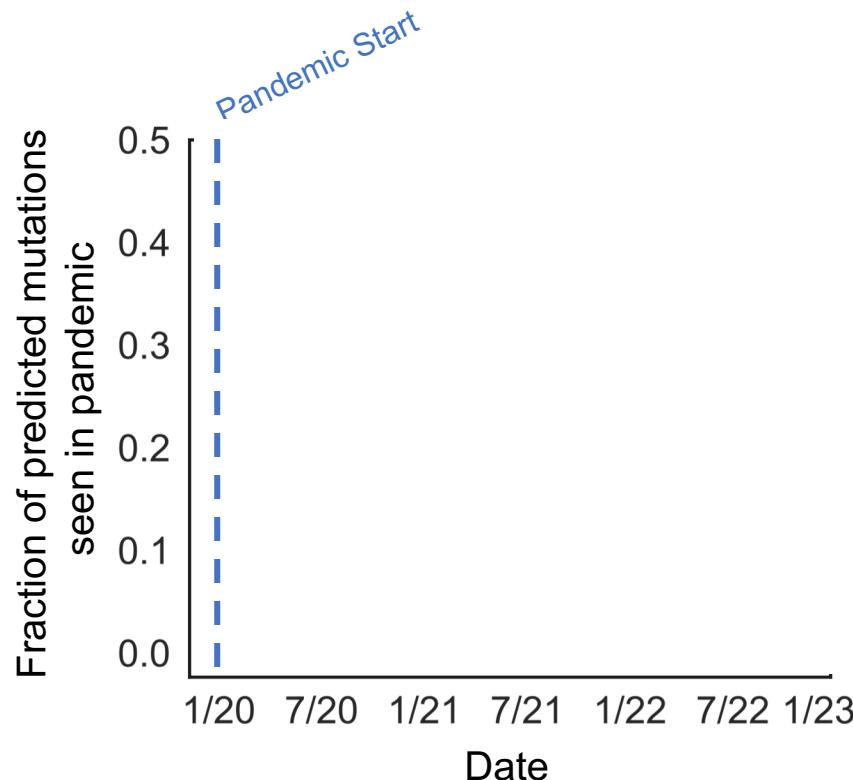
>15 million SARS-CoV-2 genomes  
>690 Spike-Ab Structures  
many experiments

Not using any data after  
Jan 2020 for model  
training

Now we can use for evaluation!

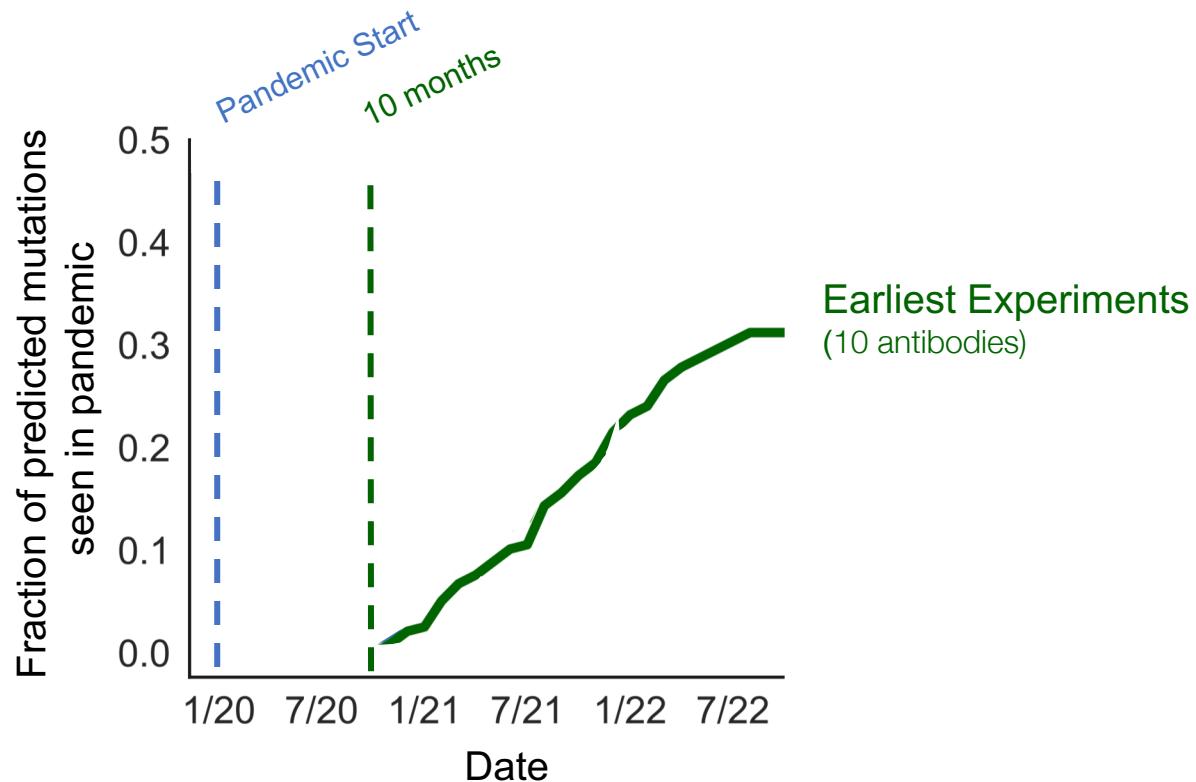
# A trip back in time to evaluate predictions

How well do current methods comprehensively predict **future immune escape?**



# A trip back in time to evaluate predictions

How well do current methods comprehensively predict **future immune escape**?



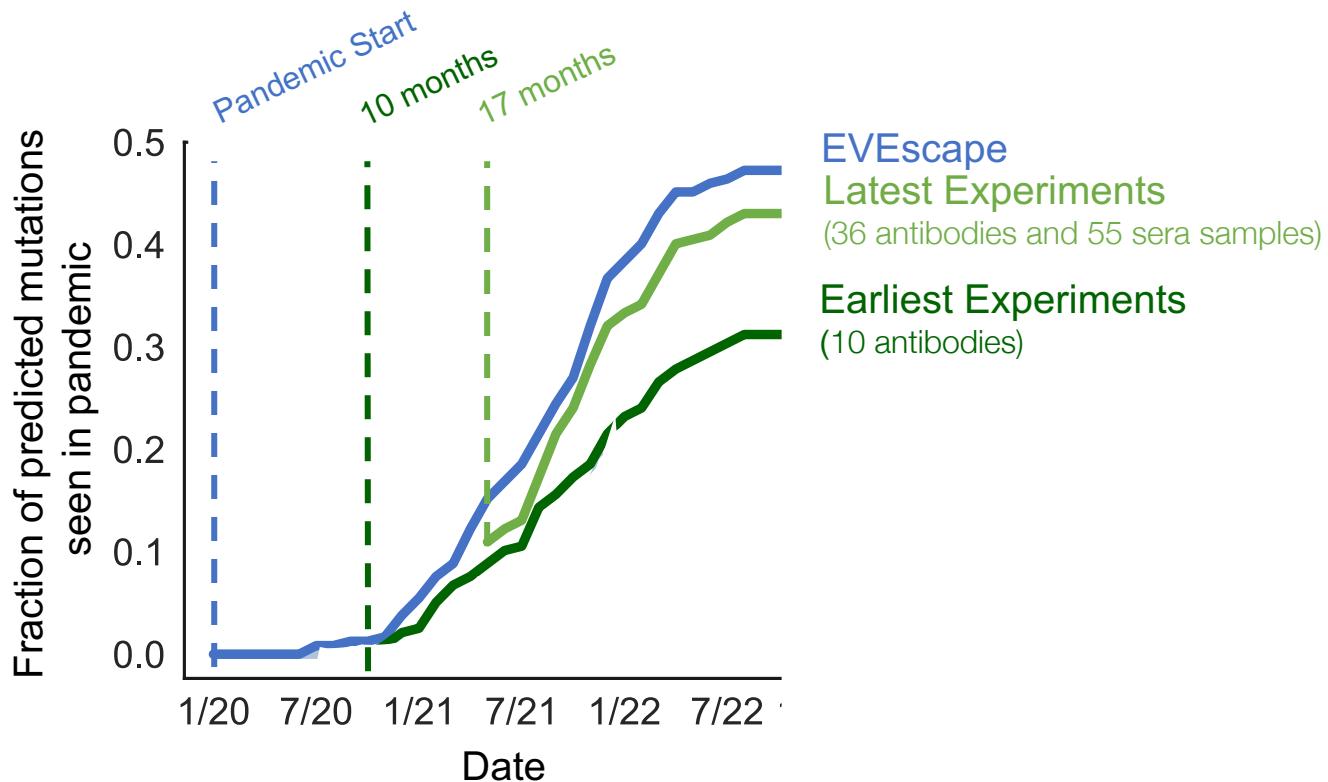
Neutralization by contemporary antibodies is a weak predictor of future immune escape mutations

# A trip back in time to evaluate predictions

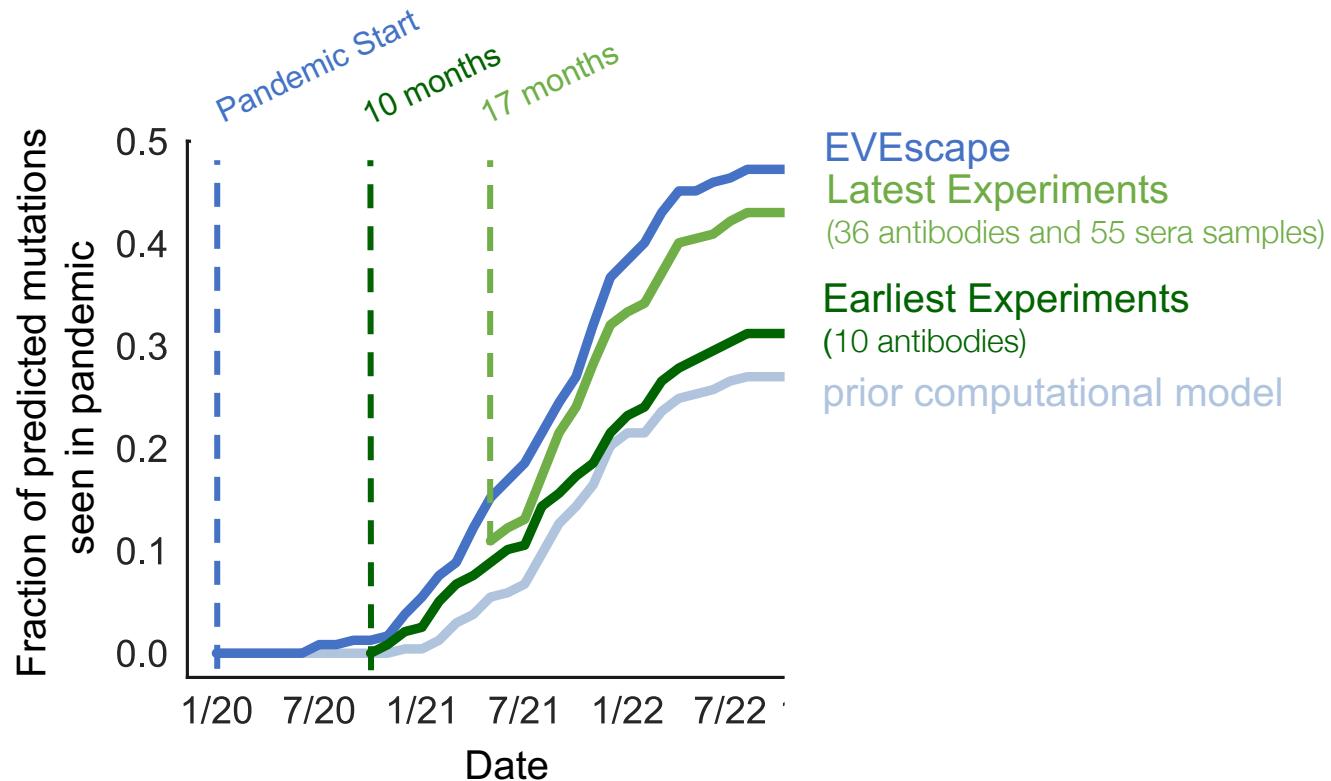
How well do current methods comprehensively predict **future immune escape**?



# EVEscape available at start of pandemic is a better predictor of viral evolution



# EVEscape available at start of pandemic is a better predictor of viral evolution



Can EVEscape trained  
on pre-pandemic data *anticipate pandemic  
mutations?*

Can EVEscape trained  
on pre-pandemic data *anticipate pandemic  
mutations?*  
Yes

# How can we apply EVEscape in current and future outbreaks?



Early warning of high-escape variants



Evaluate future protection of therapeutics



Vaccine design

# How can we apply EVEscape in current and future outbreaks?



Early warning of high-escape variants

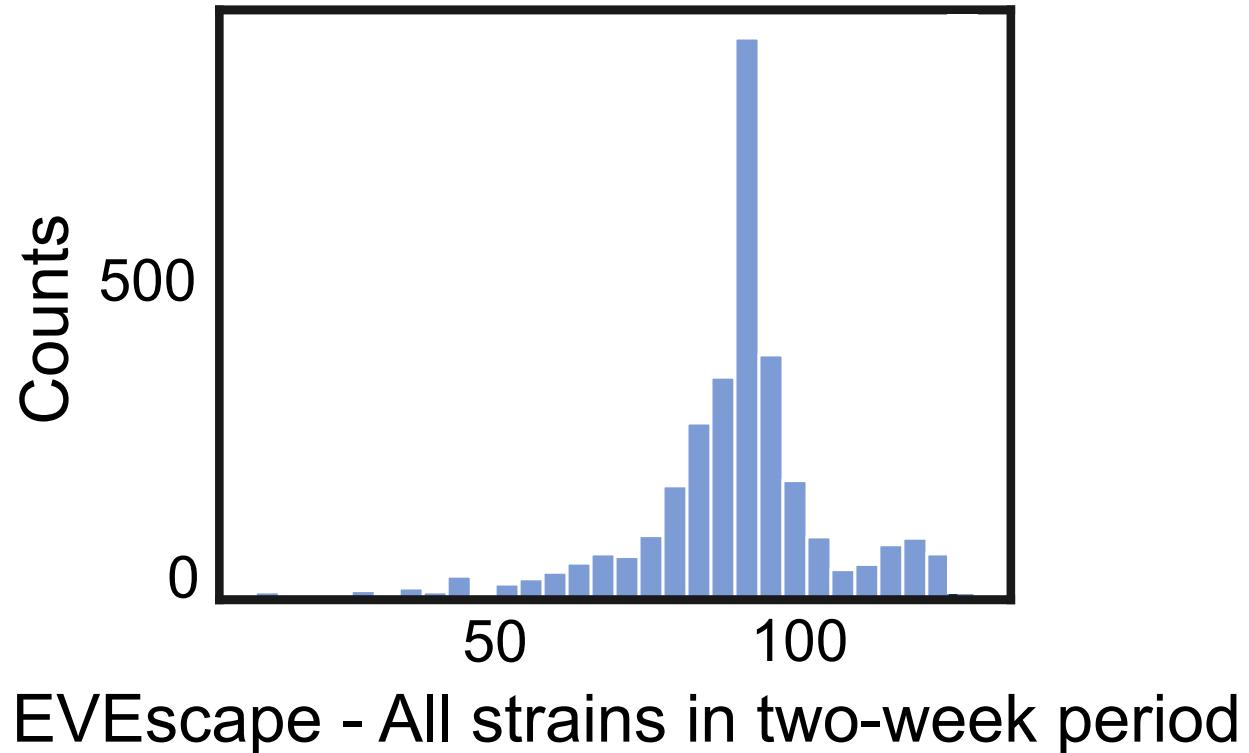


Evaluate future protection of therapeutics



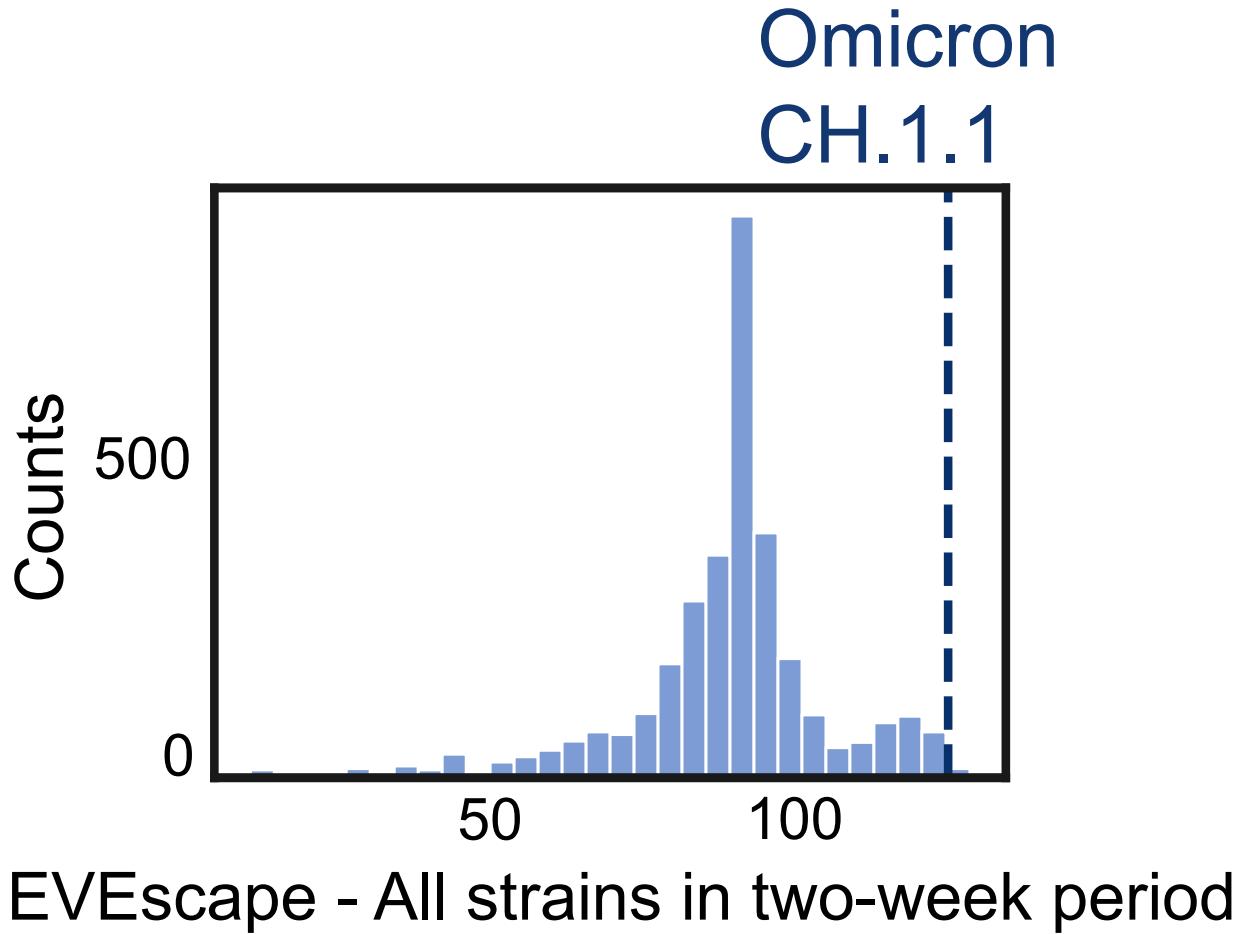
Vaccine design

Thousands of emerging strains are too many to experimentally test



# EVEscape identifies Variants of Concern out of thousands of emerging strains

[Evescape.org](https://Evescape.org)



# How can we apply EVEscape in current and future outbreaks?



Early warning of high-escape variants



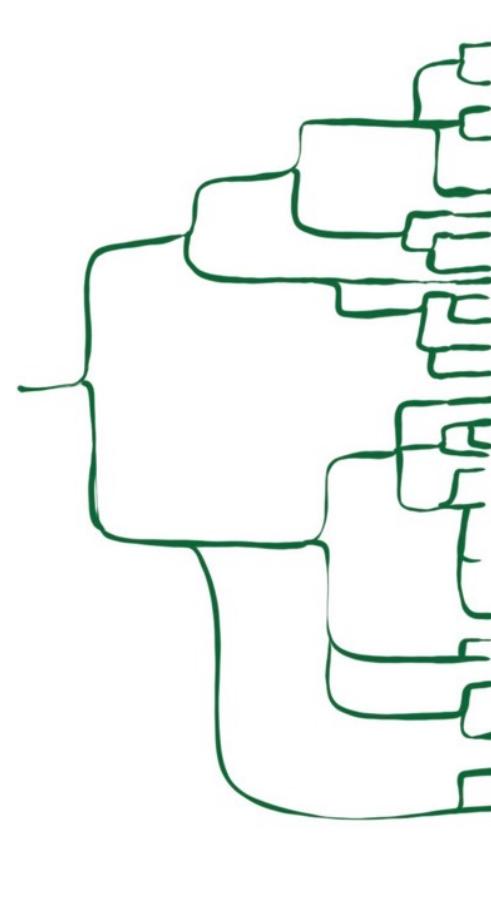
Evaluate future protection of therapeutics



Vaccine design

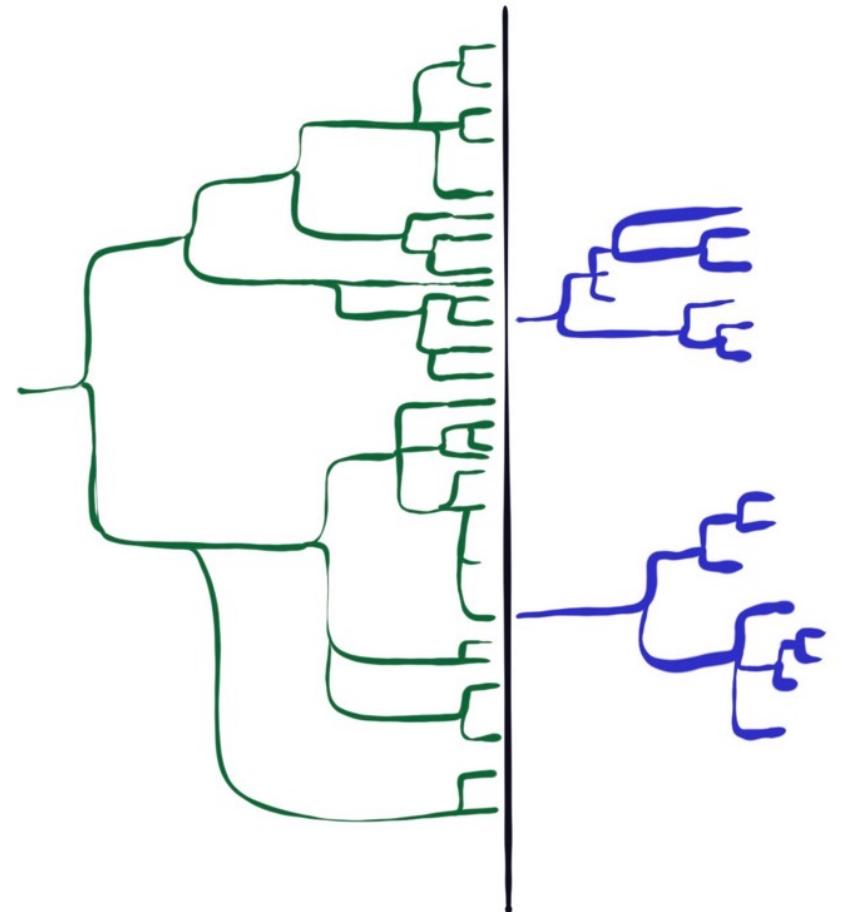
Assess vaccine efficacy  
against **existing**  
**variants**

Existing variants



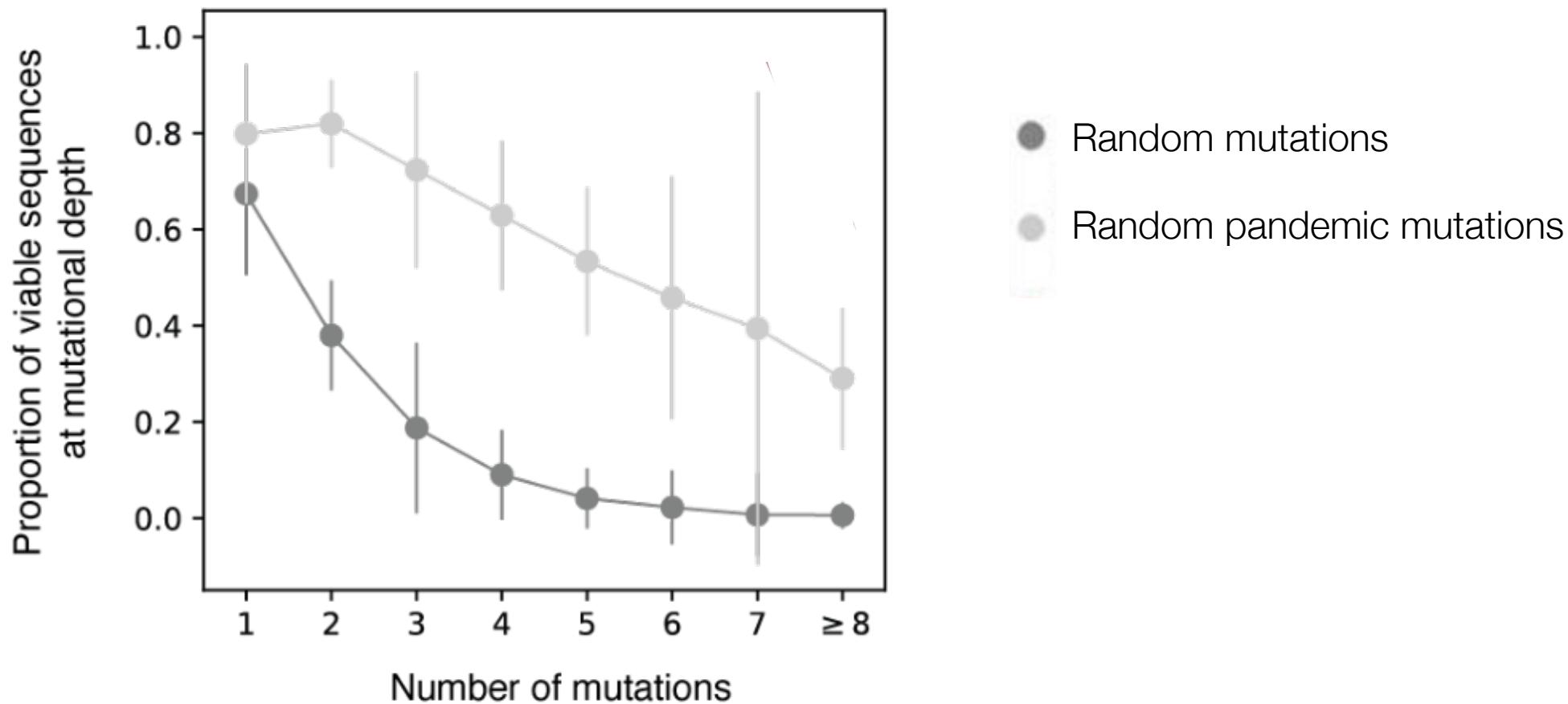
Assess vaccine efficacy  
against **existing**  
**variants**  
but also against  
**potential future variants**

Existing variants      Future variants

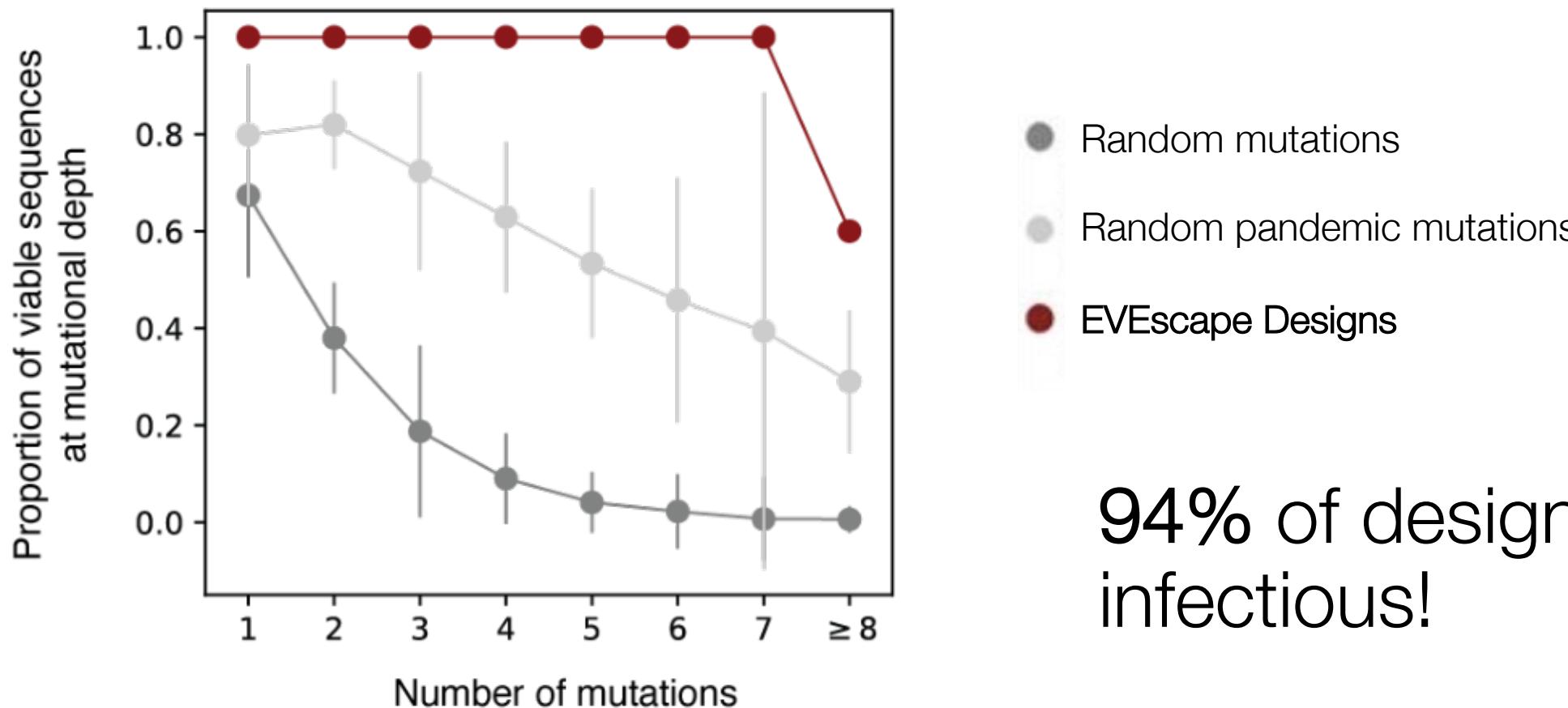


Why has this been difficult?

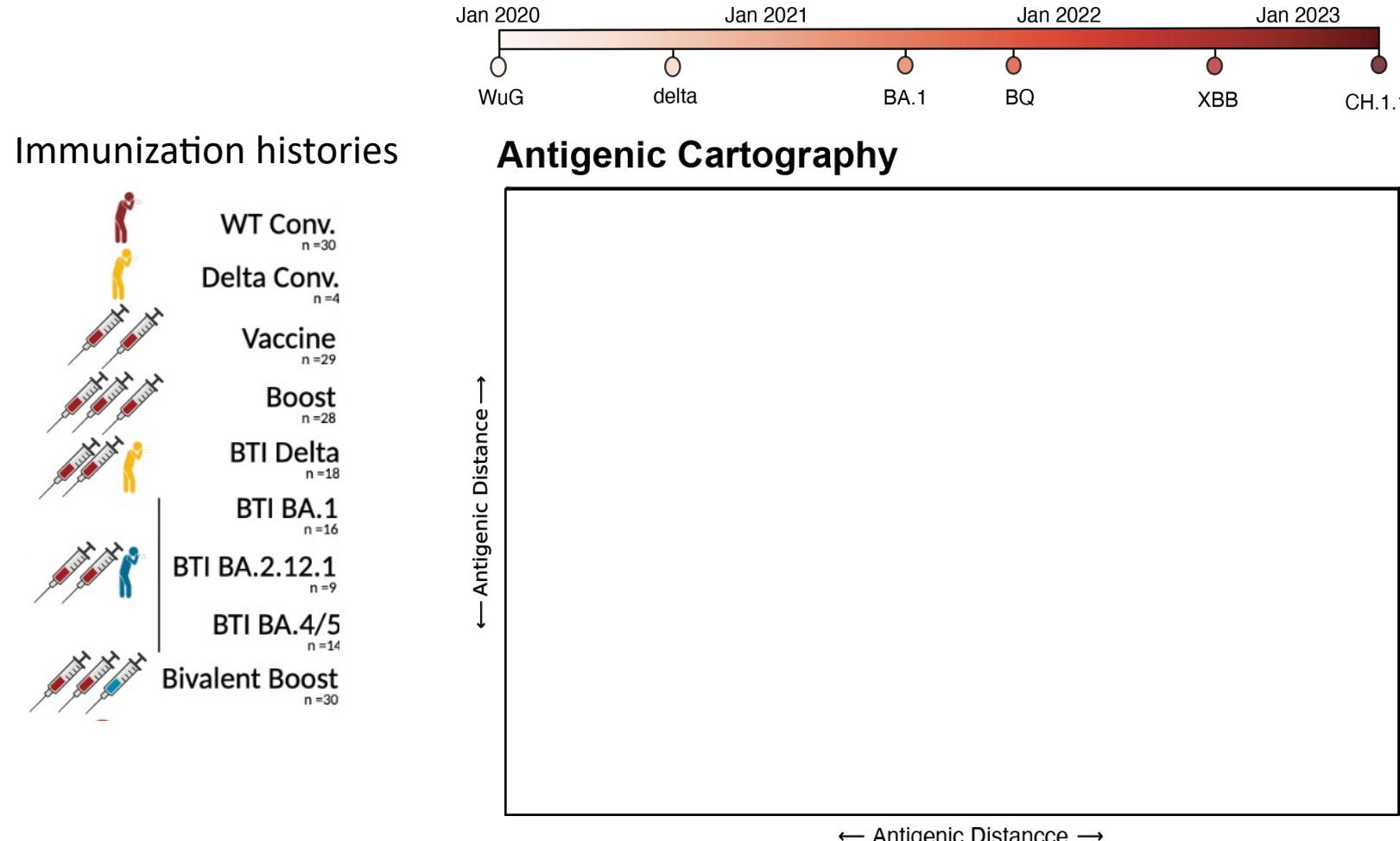
# Most random multi-mutant proteins are non-viable



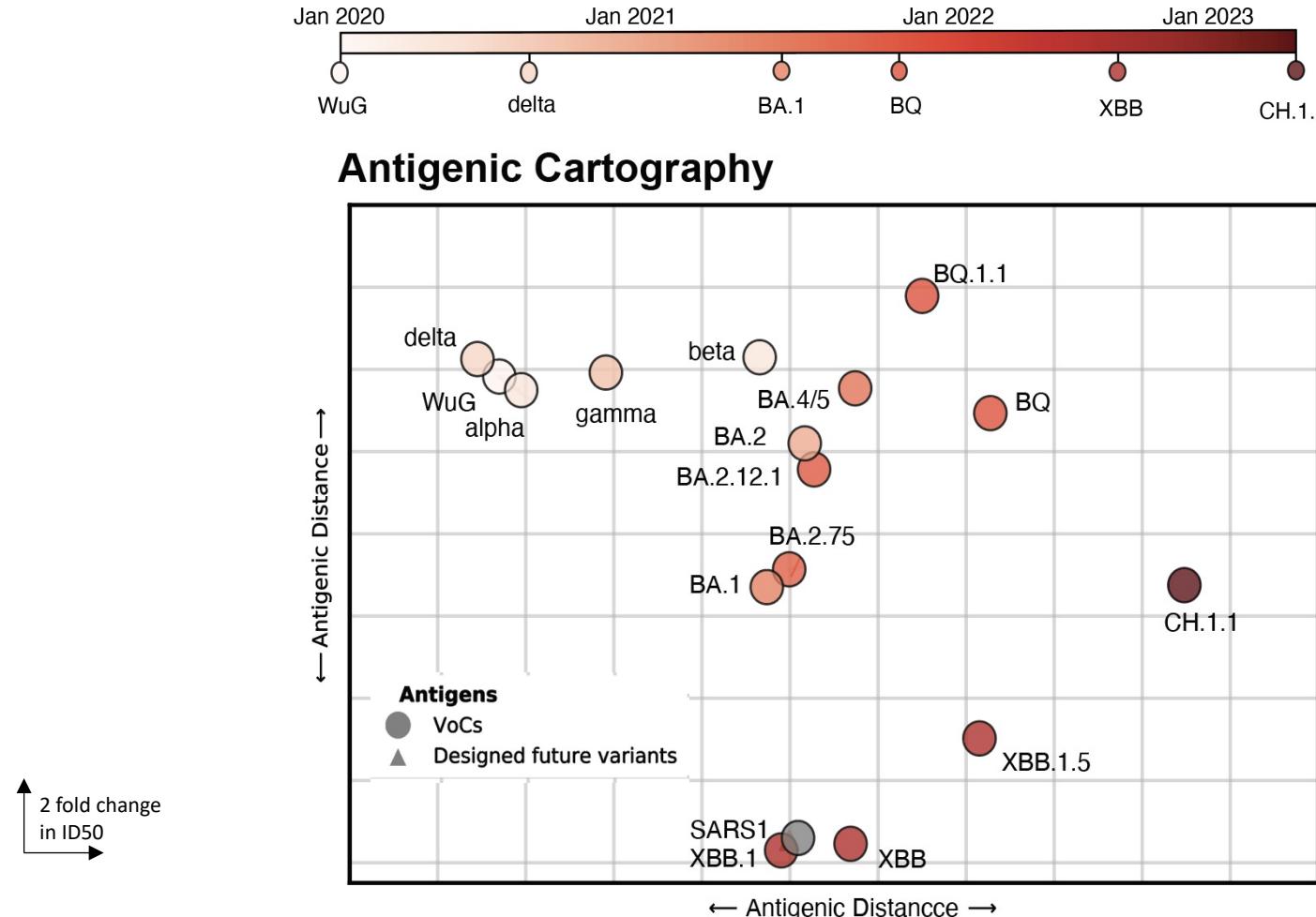
# Designing diverse multi-mutant Spikes using EVEscape



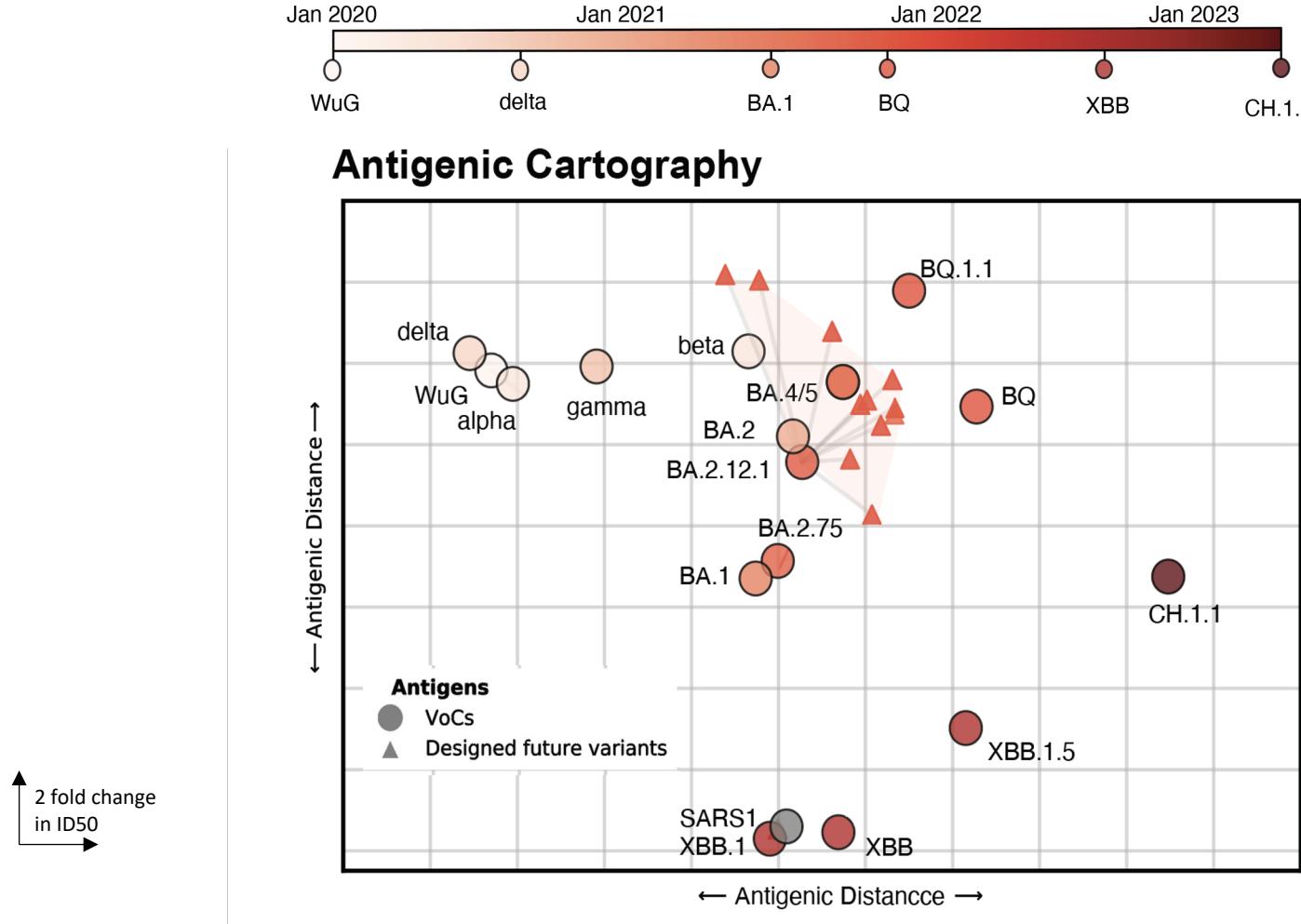
# Designed constructs mimic future evolution



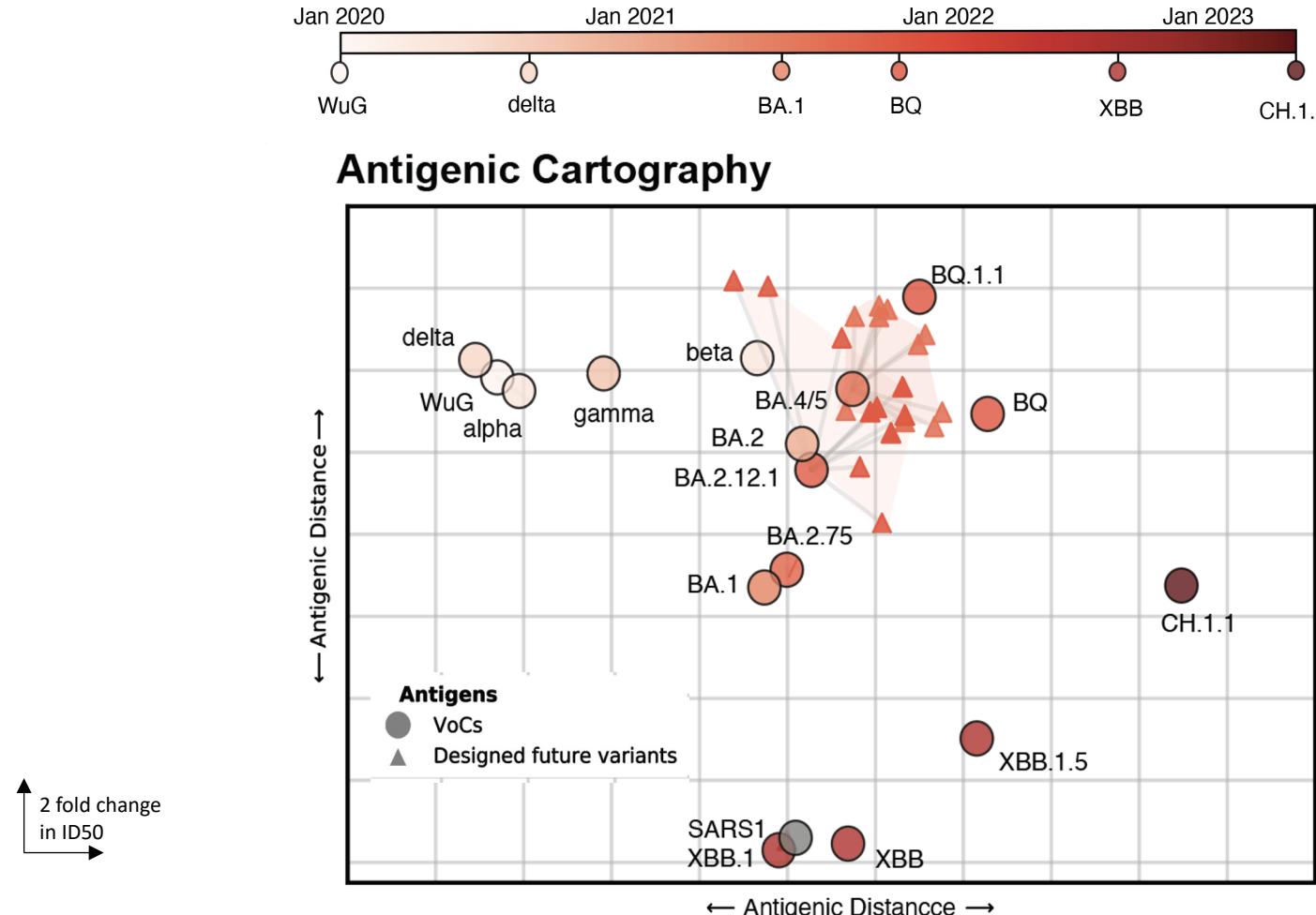
# Designed constructs mimic future evolution



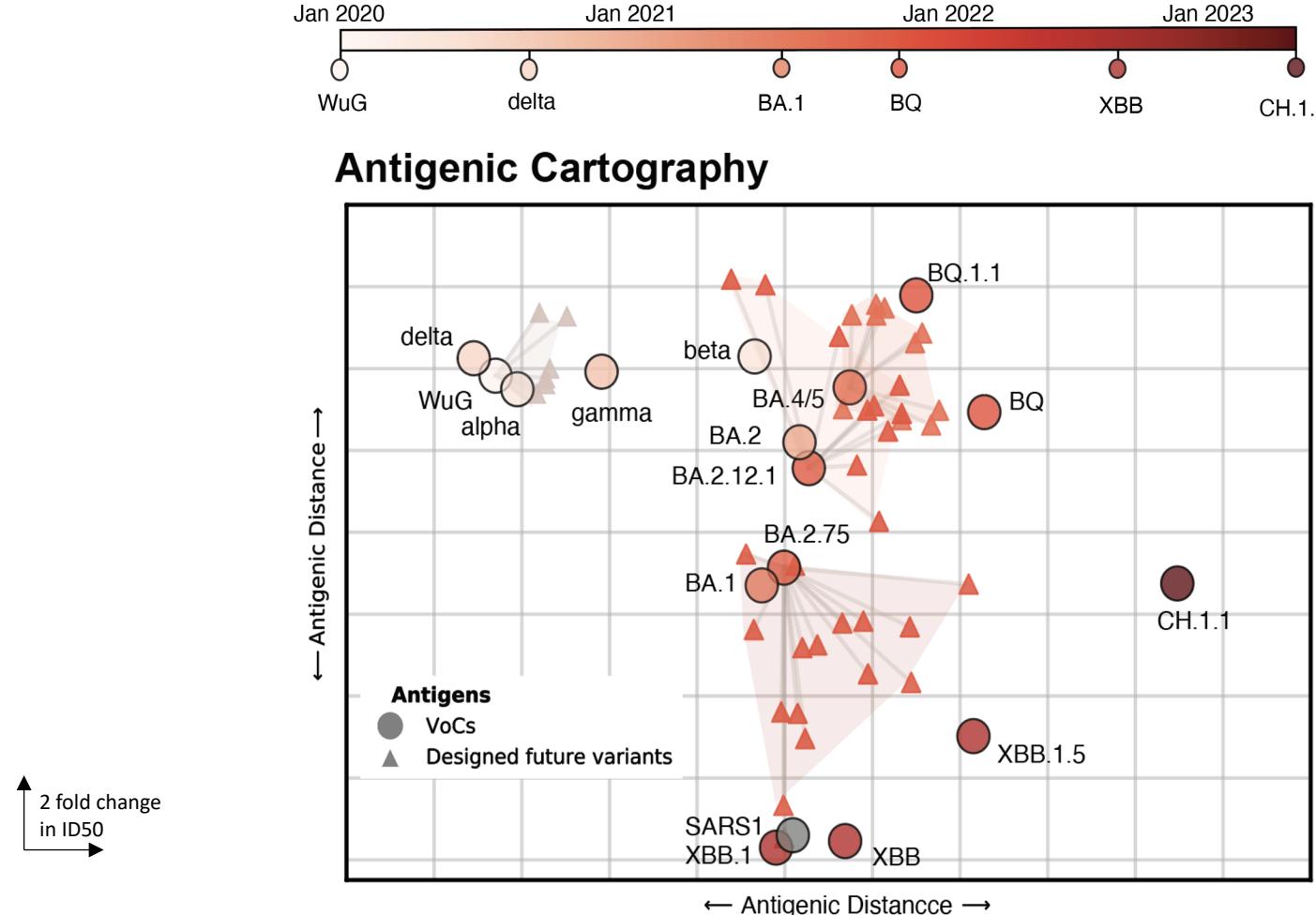
# Designed constructs mimic future evolution



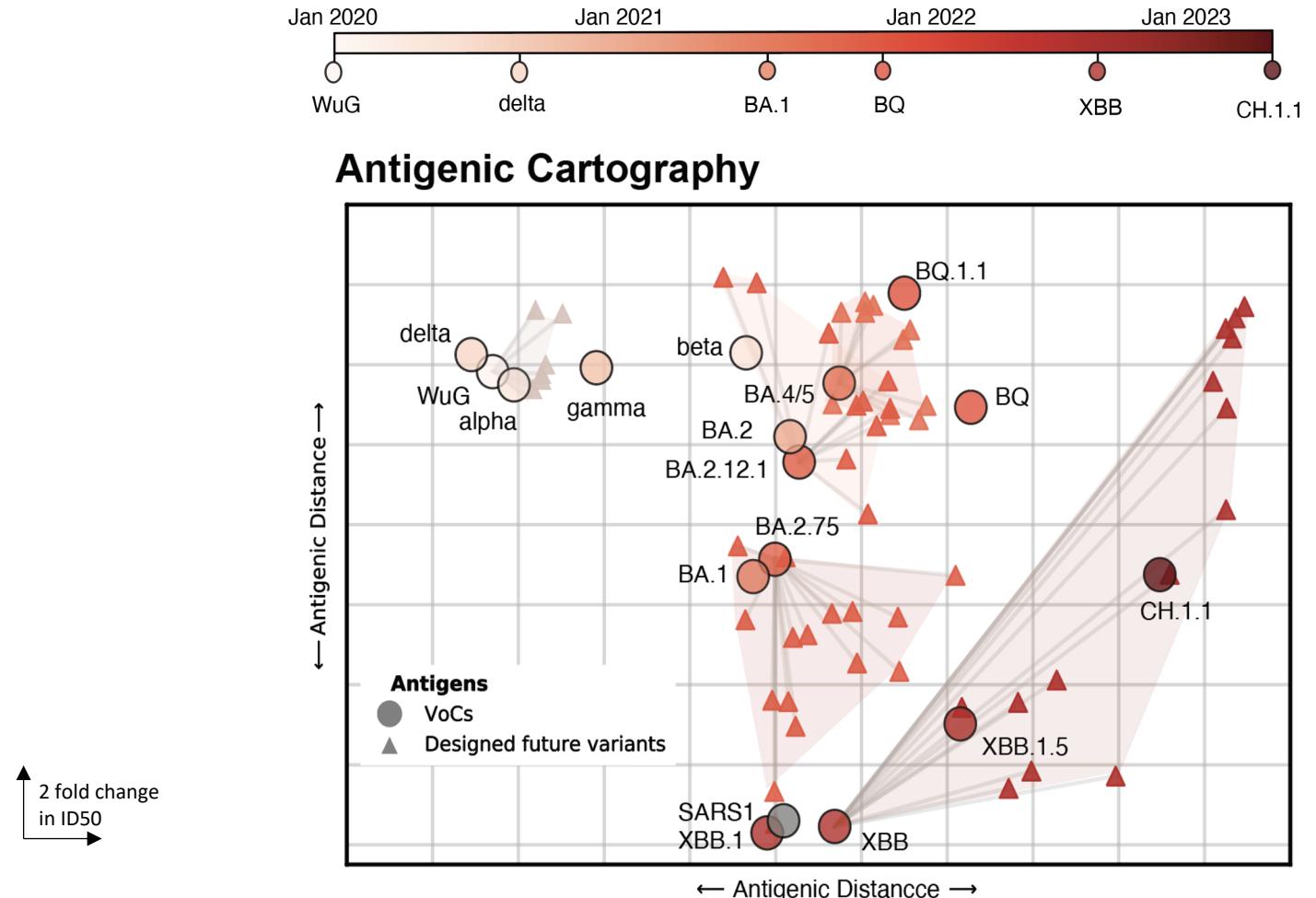
# Designed constructs mimic future evolution



# Designed constructs mimic future evolution



# Designed constructs mimic future evolution



# How can we apply EVEscape in current and future outbreaks?



Early warning of high-escape variants



Evaluate future protection of therapeutics



Vaccine design

# Computationally designed vaccines to focus antibodies to regions unlikely to mutate in the future

