

CSE7850/CX4803 Machine Learning in Computational Biology

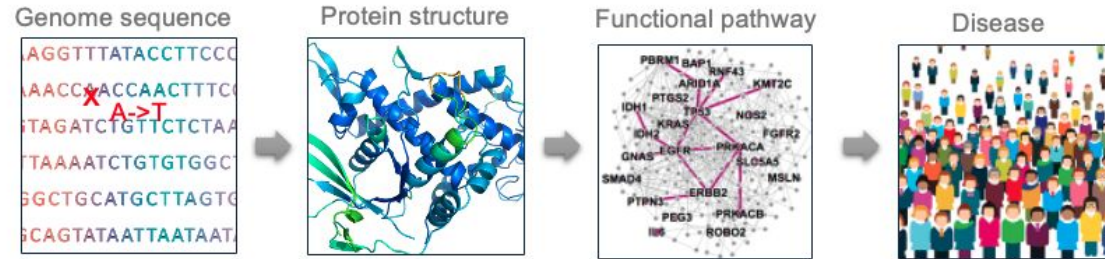


Lecture 17: Variant Effect Prediction

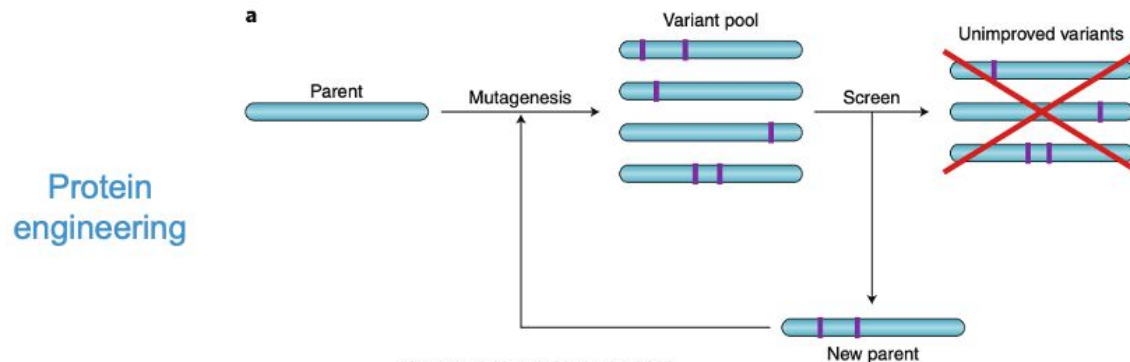
Yunan Luo

Understanding the effect of mutations/variants

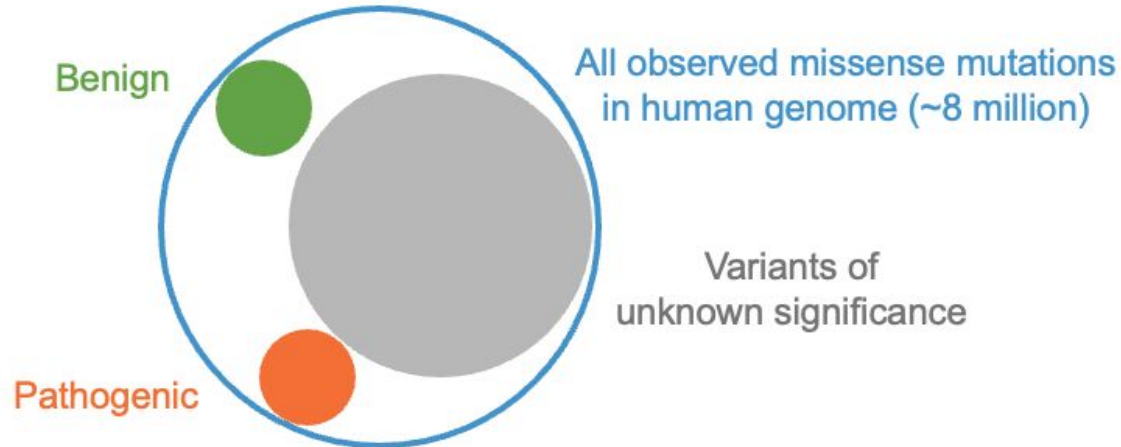
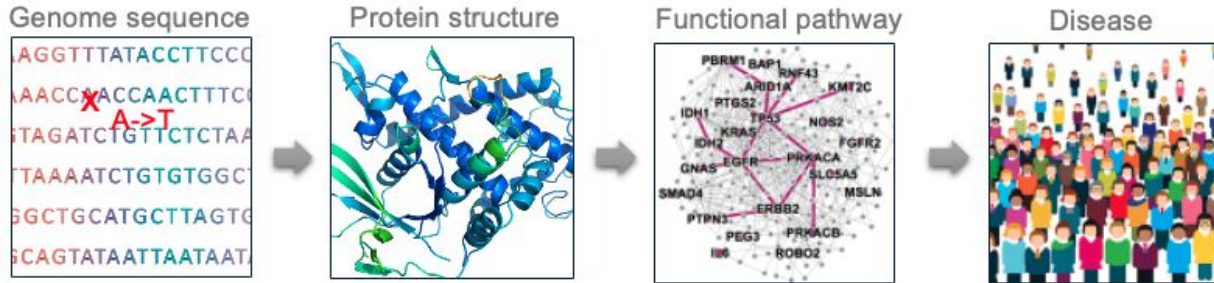
- Identify “bad” mutations (pathogenic variants in human disease)



- Identify “good” mutations (function-enhancing variants in protein engineering)

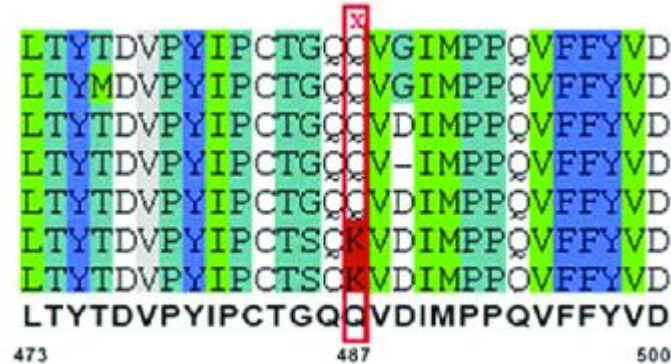


Pathogenicity of disease variants



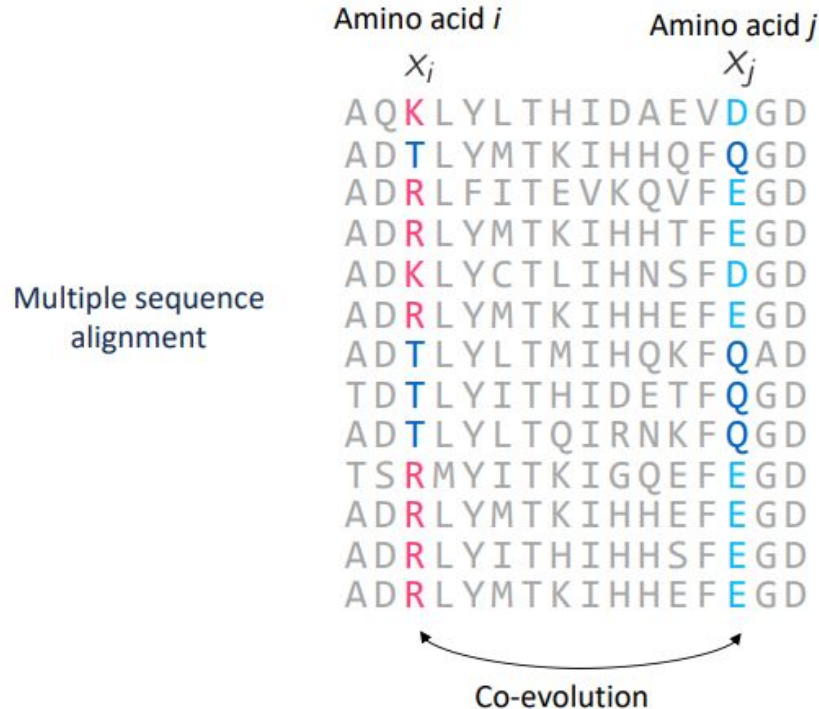
How to predict pathogenicity without labeled data?

- Evolutionary information revealed by sequence alignment:



- Any limitation of this method?

How to capture pairwise dependencies?



$$P(\mathbf{x}) = \frac{1}{Z} \exp \left(\sum_i e_i(x_i) + \sum_{i \neq j} e_{ij}(x_i, x_j) \right)$$

Single
potentials

Pairwise
potentials

Local preference

Co-evolution strength

Markov random field
Ising (Potts) model
Undirected graphical model

How to capture more than 2-order dependencies

$$P(\mathbf{x}) = \frac{1}{Z} \exp \left(\sum_i e_i(x_i) + \sum_{i \neq j} e_{ij}(x_{ij}) \right)$$

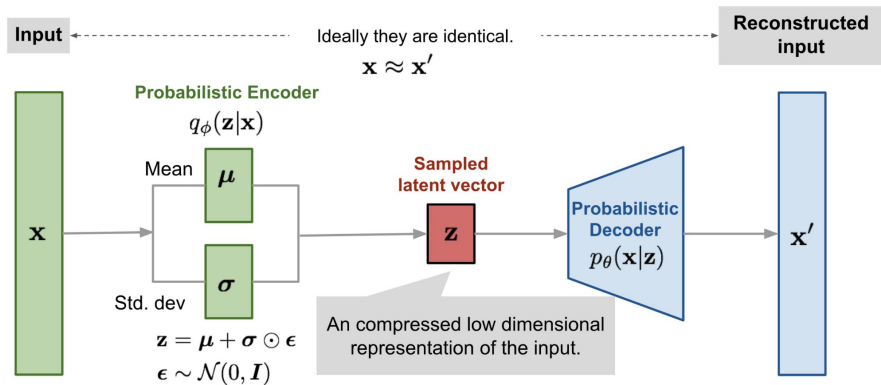
- Co-evolution based methods only consider second order interaction
 - Can we extend to model higher-order terms?

$$P(\mathbf{x}) = \frac{1}{Z} \exp \left(\sum_i e_i(x_i) + \sum_{i \neq j} e_{ij}(x_{ij}) + \sum_{i \neq j \neq k} e_{ijk}(x_{ijk}) + \dots \right)$$

- Computationally hard
 - An exponential number of terms
- Latent-variable ML models capture higher-order interactions *implicitly*

$$p(\mathbf{x}|\boldsymbol{\theta}) = \int p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z}$$

Pseudocode of VAE



```

1 class variational_autoencoder(nn.Module):
2     def __init__(self, ...):
3         super(variational_autoencoder, self).__init__()
4         self.encoder = Encoder(...)
5         self.decoder = Decoder(...)
6
7     def forward(self, x):
8         mu, sigma = self.encoder(x)
9         z = self.sample(mu, sigma)
10        x_prime = self.decoder(z)
11        return x

```

- Objective function: $\max \mathbb{E}_q [\log p(x|z)] - D_{\text{KL}}(q||p)$
- Reconstruction consistency:** $\mathbb{E}_q[\log p(x|z)] = -\frac{1}{2\sigma^2} \mathbb{E}_q[\|x - G_\theta(z)\|^2] + \text{const}$
- KL divergence between $p(z)$ and $q(z)$:** typically, $p(z) = \mathcal{N}(0, 1)$; The KL term encourages $q(z)$ to be close to the standard normal distribution $\mathcal{N}(0, 1)$

Paper #2

Article

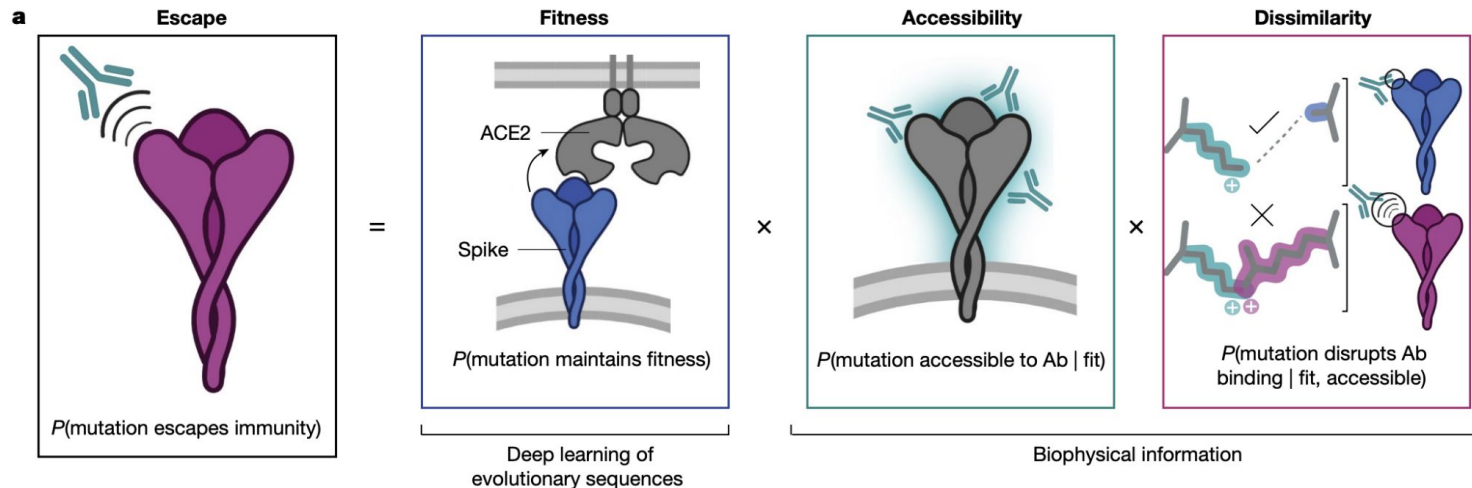
Learning from prepandemic data to forecast viral escape

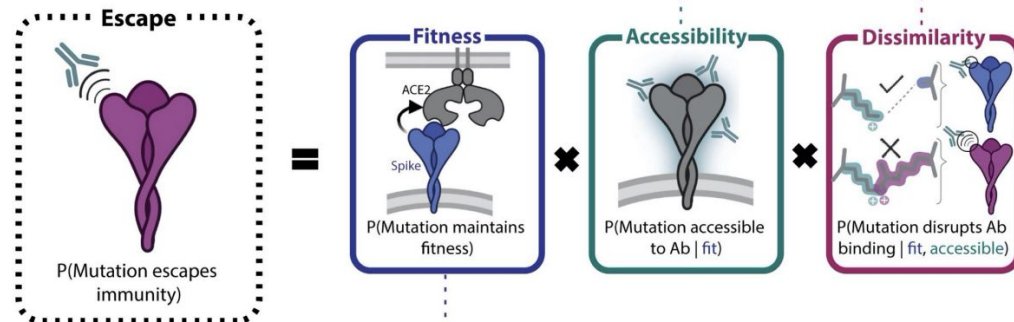
<https://doi.org/10.1038/s41586-023-06617-0>

Received: 20 July 2022

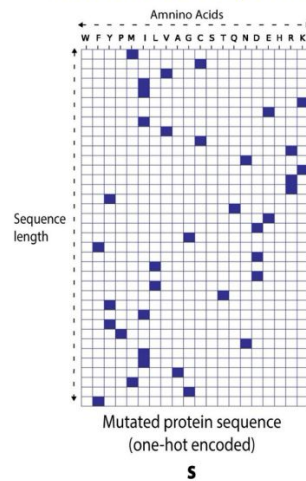
Accepted: 6 September 2023

Nicole N. Thadani^{1,6}, Sarah Gurev^{1,2,6}, Pascal Notin^{3,6}, Noor Youssef¹, Nathan J. Rollins^{1,5}, Daniel Ritter¹, Chris Sander^{1,4}, Yarin Gal³ & Debora S. Marks^{1,4}✉

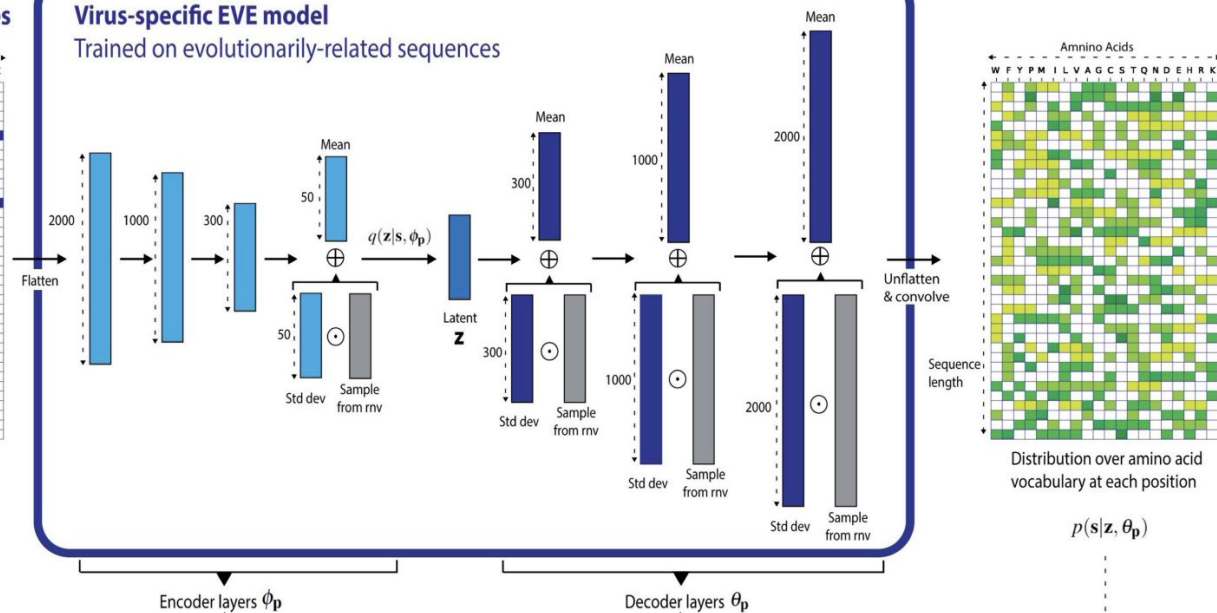




Fitness component input:
Mutated & WT sequences



Virus-specific EVE model
Trained on evolutionarily-related sequences



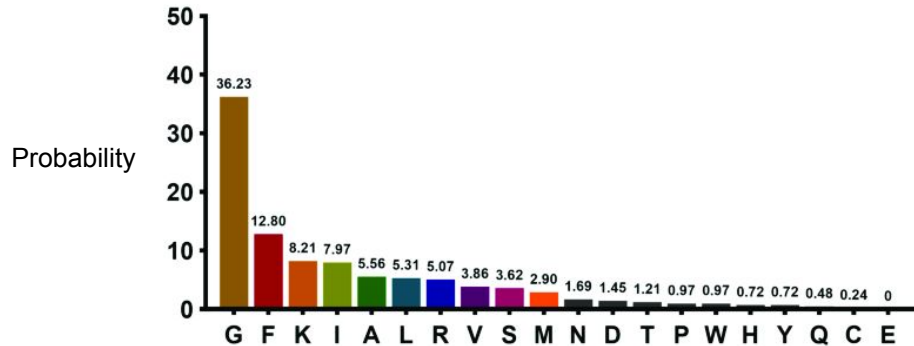
Another way to predict mutation effect: language model

Masked language models

$$p(x) = \prod_{i=1}^L p(x_i | x_1 \dots x_{i-1}, x_{i+1} \dots x_L)$$

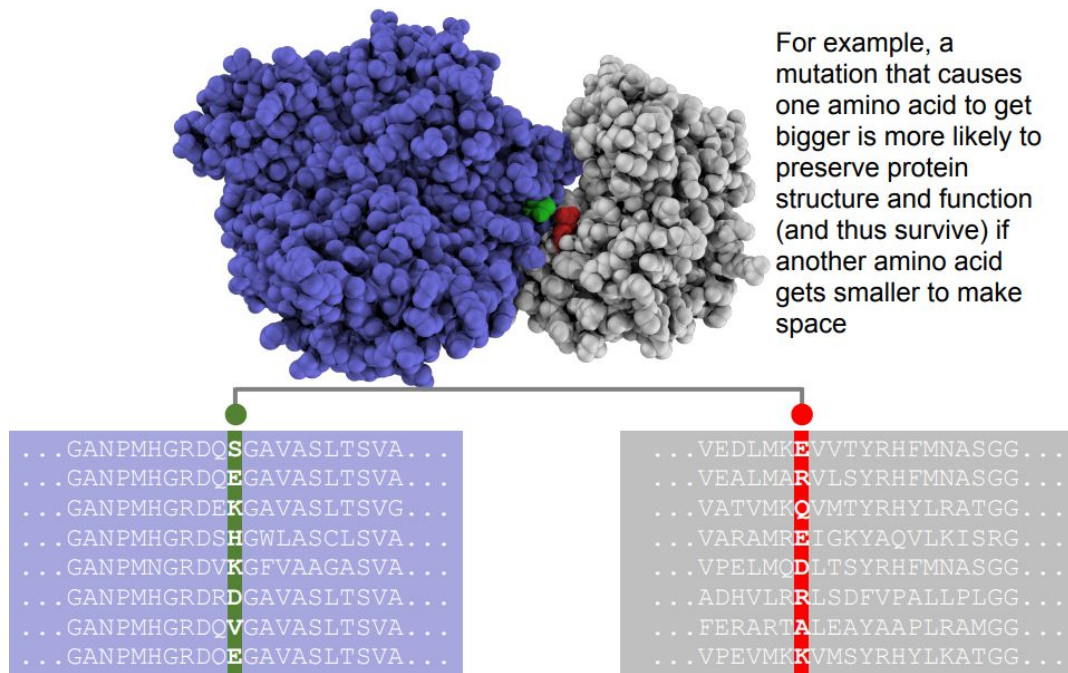
- Language models capture higher-order interactions *implicitly*

MSKGEE??TGVVPI????DGDVNGHKFSVY



PLMs are trained on natural sequence data ...

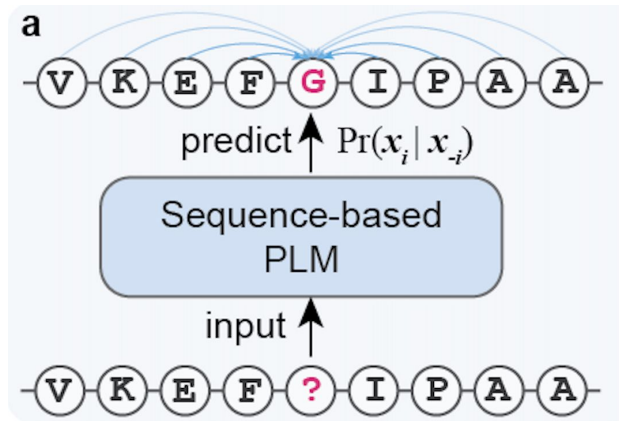
but we have seen that structure can have an “evolutionary effect” on protein sequences



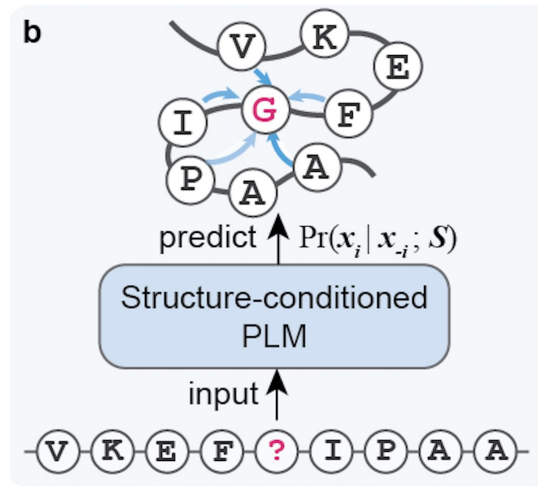
Can we leverage structure data in protein language models?

Structure-based PLM

Traditional PLM (sequence-based)



Structure-based PLM



Paper #1

RESEARCH

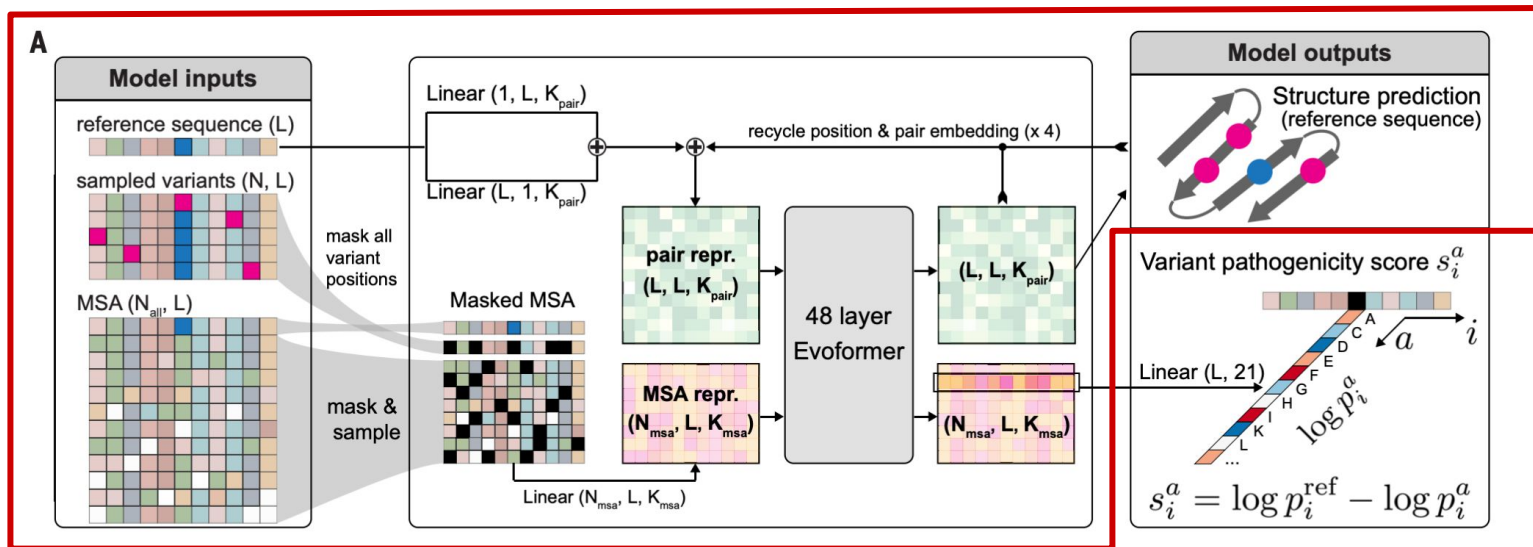
RESEARCH ARTICLE SUMMARY

MACHINE LEARNING

Accurate proteome-wide missense variant effect prediction with AlphaMissense

Jun Cheng*, Guido Novati, Joshua Pan†, Clare Bycroft†, Akvilė Žemgulytė†, Taylor Applebaum†, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hassabis, Pushmeet Kohli*, Žiga Avsec*

Almost identical to AlphaFold



Lecture:

Topic	Contents
Introduction Basics in computational biology	Introduction & Logistics
	Molecular biology
	No class (MLK day)
	Sequence alignment I
	Sequence alignment II
ML foundations	No Class (PyTorch video + exercise)
	Regression & Gradient descent
	Classification & Toolbox for Applied ML
	Neural networks
	Deep learning
Learning from sequence data	Deep learning for Protein/DNA sequences
	Large language models (LLMs)
Learning from high-dim data	Clustering and dimensionality reduction
	Generative AI
Learning from network data	Network basics & ML for graphs
	Graph neural network
Learning from structure data	Protein structure prediction & protein design
Advanced topics: ML for sequence data	Protein language models for prediction and generation
	Disease variant prediction
	Protein function prediction
	No class (Spring break)
	No class (Spring break)
Advanced topics: ML for structure data	Deep learning for structure prediction
	GNN for 3D structures
	Deep learning for structure generation
Advanced topics: ML for network data	Embeddings (representation learning)
	ML for protein design
	ML for drug discovery
Advanced topics: ML for high-D data	Dim reduction in bio data
	ML for system bio
	ML-guided biological discovery

Biology background

Intro to ML and DL

ML for bio data

ML in CompBio research

Hands-on exercise:

 PyTorch

 Google colab

 kaggle

Announcements: Midterm Survey

- Informal Midterm Course Feedback (anonymous)
 - Provide mid-term feedback for this course
 - Suggest your favorite topics! May be incorporated in the remaining lectures
- Access the survey form by one of the following:
 - <https://tinyurl.com/mlb-ief-s24>
 - Canvas -> Syllabus -> “Midterm Survey”
 - QR Code:

