

BIOCHEMISTRY & MOLECULAR BIOLOGY

Seminar Series

Scalable Emulation of Protein Equilibrium Ensembles *with* Generative Deep Learning

ANDREW Y. K. FOONG, PH.D.

October 14th, 2025

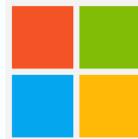


Radiation
Oncology
AI & Data Analytics
AIDA

About me



Senior Associate Consultant · Radiation Oncology



Senior Researcher · Microsoft Research



Ph.D. in Machine Learning · Cambridge University

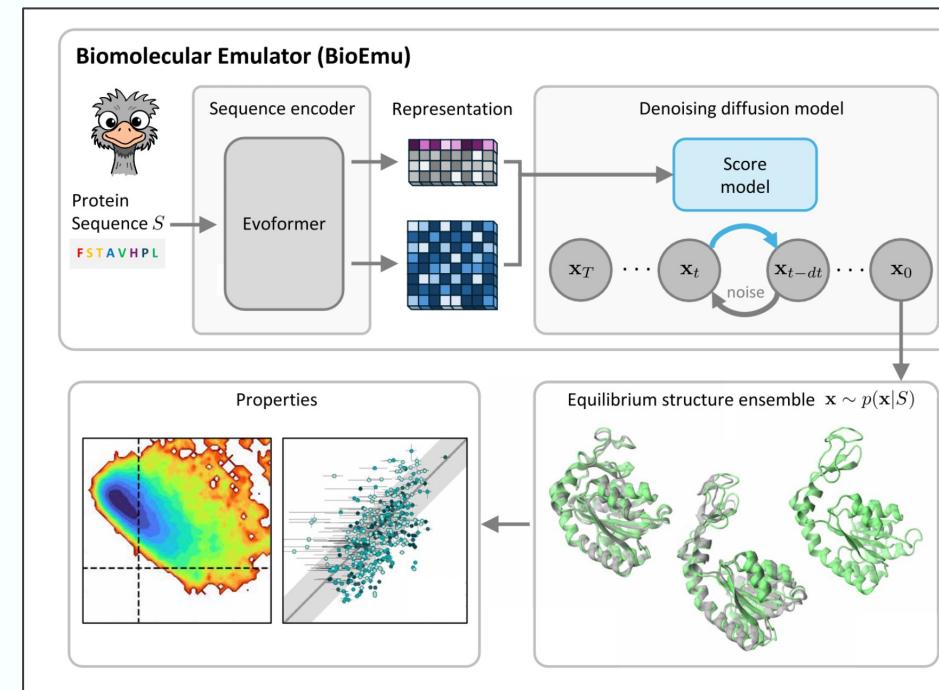


Research Scientist Intern · Google DeepMind

More info: <https://andrewfoongyk.github.io>

BioEmu—a Biomolecular Emulator

- BIOEMU is a deep learning model that takes protein sequences as input, and outputs a *distribution* over 3D protein structures.
- The distributions are trained to emulate physical conformational variability:
 - Predicts *multiple* conformations
 - Predicts *thermodynamics*: the relative probabilities of each conformation.



BioEmu—a Biomolecular Emulator



PROTEIN SIMULATIONS

Scalable emulation of protein equilibrium ensembles with generative deep learning

Sarah Lewis^{1†}, Tim Hempel^{1†}, José Jiménez-Luna^{1†}, Michael Gastegger^{1†}, Yu Xie^{1†}, Andrew Y. K. Foong^{1†}, Victor García Satorras^{1†}, Osama Abdin^{1†}, Bastiaan S. Veeling^{1†}, Iryna Zaporozhets^{1,2}, Yaoyi Chen^{1,2}, Soojung Yang¹, Adam E. Foster¹, Arne Schneuing¹, Jigyasa Nigam¹, Federico Barbero¹, Vincent Stumper¹, Andrew Campbell¹, Jason Yim¹, Marten Lienen¹, Yu Shi¹, Shuxin Zheng¹, Hannes Schulz¹, Usman Munir¹, Roberto Sordillo¹, Ryota Tomioka¹, Cecilia Clementi^{1,2,3}, Frank Noé^{1,2,3*}

Following the sequence and structure revolutions, predicting functionally relevant protein structure changes at scale remains an outstanding challenge. We introduce BioEmu, a deep learning system that emulates protein equilibrium ensembles by generating thousands of statistically independent structures per hour on a single graphics processing unit (GPU). BioEmu integrates more than 200 milliseconds of molecular dynamics (MD) simulations, static structures, and experimental protein stabilities using new training algorithms. It captures diverse functional motions—including cryptic pocket formation, local unfolding, and domain rearrangements—and predicts relative free energies with 1 kilocalorie per mole accuracy compared with millisecond-scale MD and experimental data. BioEmu provides mechanistic insights by jointly modeling structural ensembles and thermodynamic properties. This approach amortizes the cost of MD and experimental data generation, demonstrating a scalable path toward understanding and designing protein function.

Sarah Lewis, Tim Hempel, José Jiménez-Luna, Michael Gastegger, Yu Xie, Andrew Y. K. Foong, Victor García Satorras, Osama Abdin, Bastiaan S. Veeling *et al.*, Scalable emulation of protein equilibrium ensembles with generative deep learning. *Science* 389, eadv9817 (2025). DOI: [10.1126/science.adv9817](https://doi.org/10.1126/science.adv9817)

Today's talk

1. Background

- *AlphaFold and the deep learning revolution*

2. AlphaFold

- *How does it work?*

3. Structure is not enough

- *Equilibrium distributions*

4. BioEmu

- *Training data*
- *Diffusion model*
- *Capabilities*

5. Q&A

Background

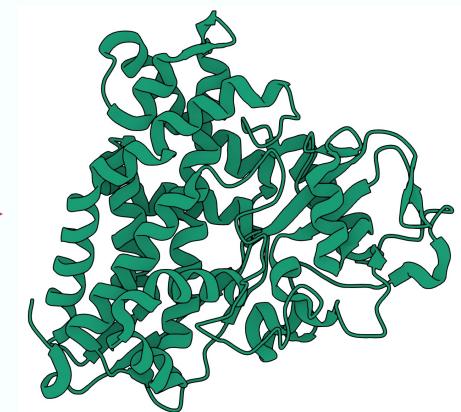
AlphaFold and the Deep Learning Revolution

The structure prediction problem

- Protein *function* is mediated through structure.
- Protein *structure* is determined by its amino acid sequence.
- PROBLEM: Given an amino acid sequence, *what is the three-dimensional structure that the protein folds into?*

NLAPLPPHVPEHLVFDFDM...

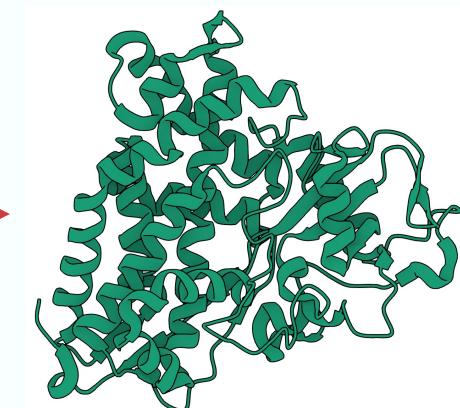
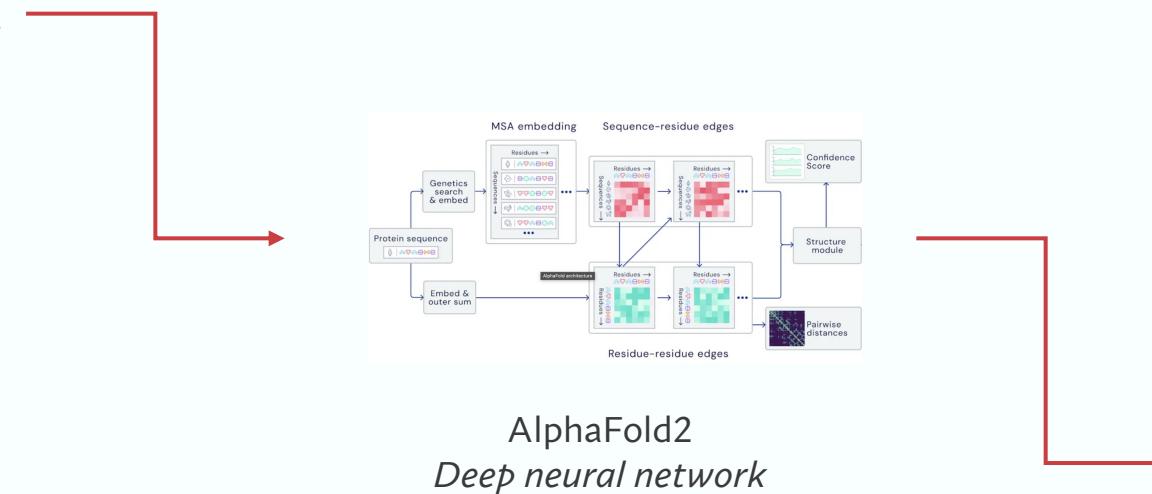
- X-ray crystallography
- Cryo-electron microscopy (cryo-EM)
- Nuclear magnetic resonance (NMR) spectroscopy



Deep learning for protein structure

- Google DeepMind cast protein structure prediction as *supervised learning*.
- Training data obtained from the Protein DataBank (PDB).
- KEY BREAKTHROUGH: *deep learning*.

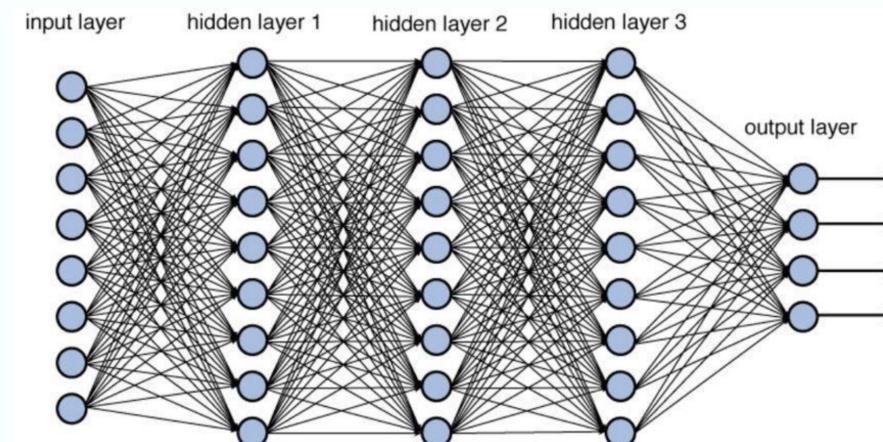
NLAPLPPHVPEHLVFDFDM...



Jumper, J., Evans, R., Pritzel, A. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021). <https://doi.org/10.1038/s41586-021-03819-2>

Deep learning in a nutshell

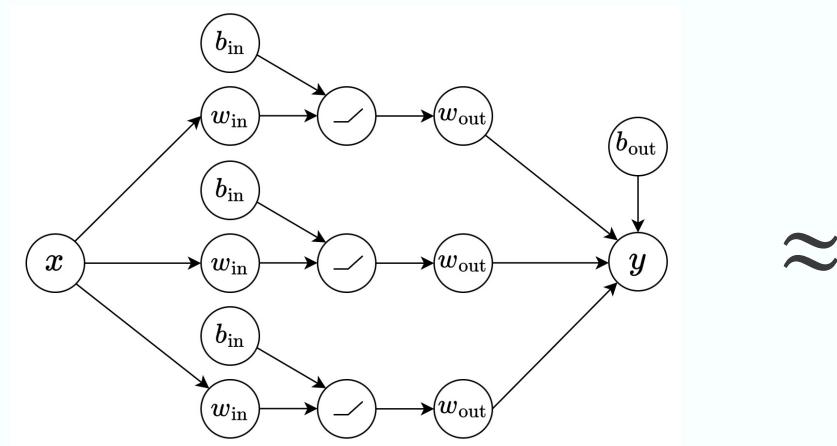
- Deep learning = use of *artificial neural networks*.
- Mathematical functions with millions of numbers: “parameters/weights”.
- The numbers determine how the neural network behaves.
 1. Start by choosing parameters randomly (garbage predictions).
 2. Optimizer automatically adjusts parameters to fit example data.
 3. Apply the function to new data (great predictions, hopefully).



Deep learning in a nutshell

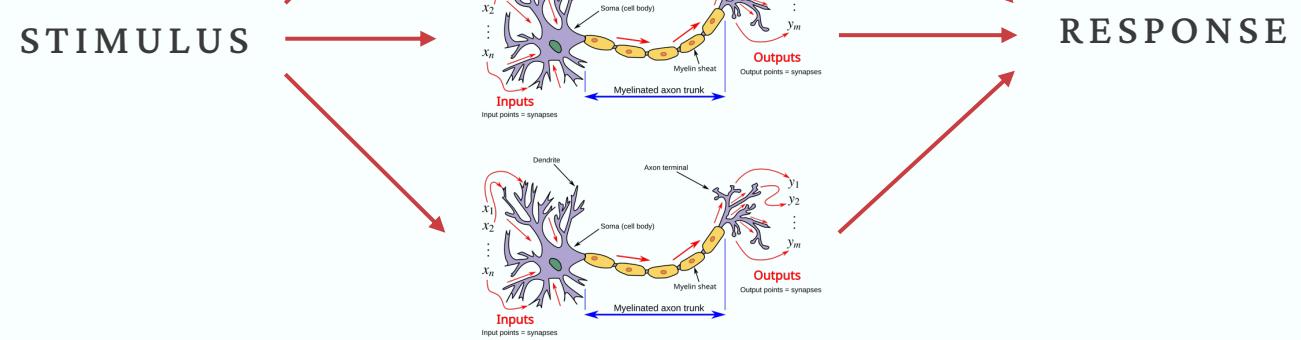
- Neural networks composed of hundreds of thousands of *artificial neurons*.
- Each neuron performs linear regression.
- Layers of these neurons are stacked interleaved with non-linear *activation functions*.

TOY EXAMPLE NETWORK (*vastly simplified*)

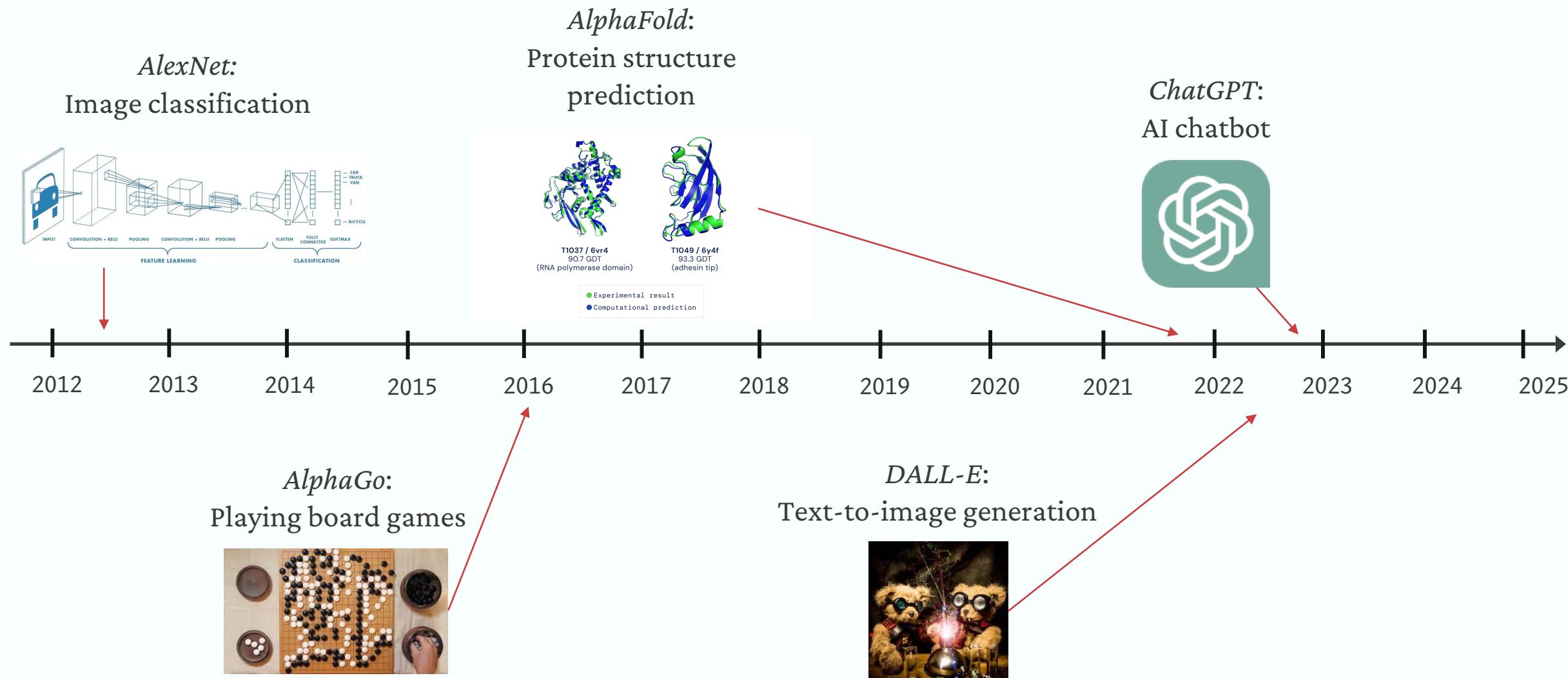


\approx

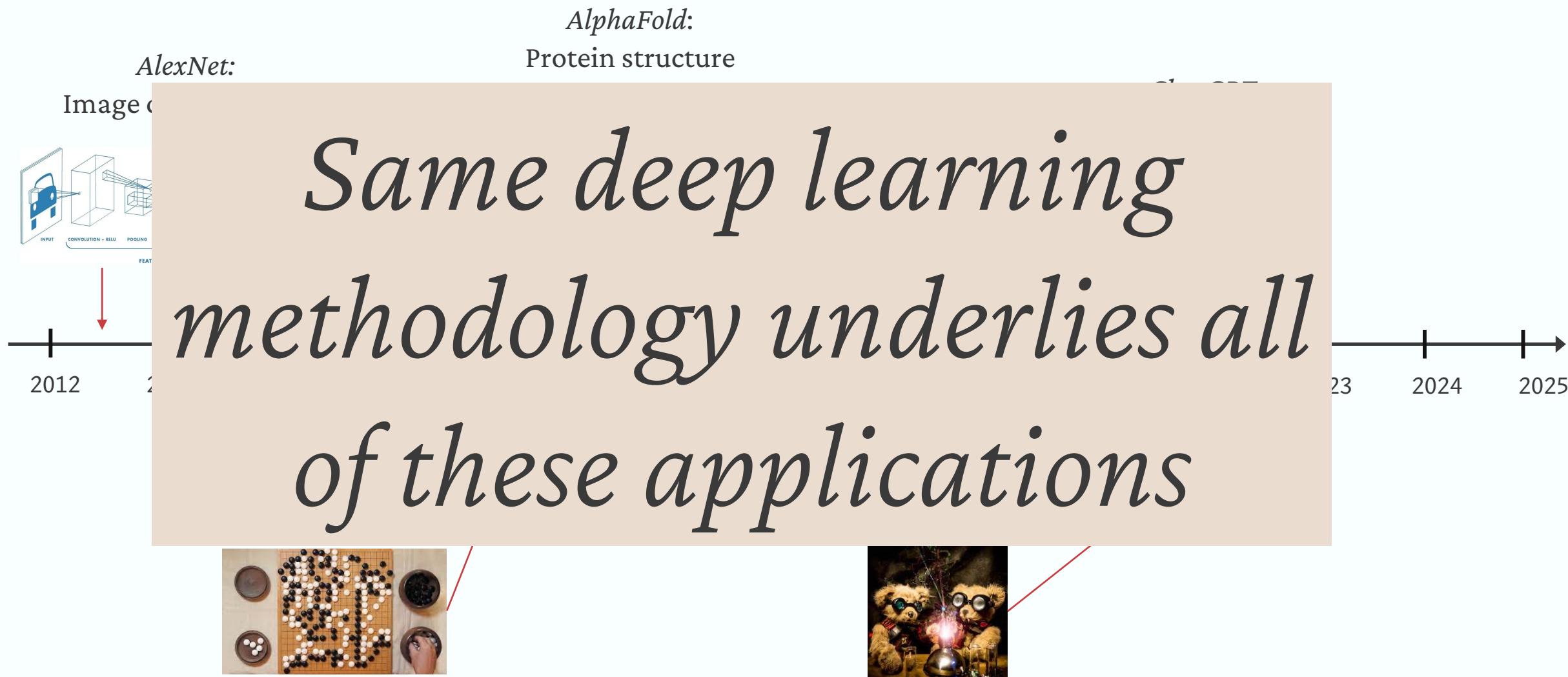
STIMULUS



Deep learning timeline



Deep learning timeline



Nobel prize for AlphaFold



DEMIS HASSABIS
Google DeepMind

JOHN JUMPER
Google DeepMind

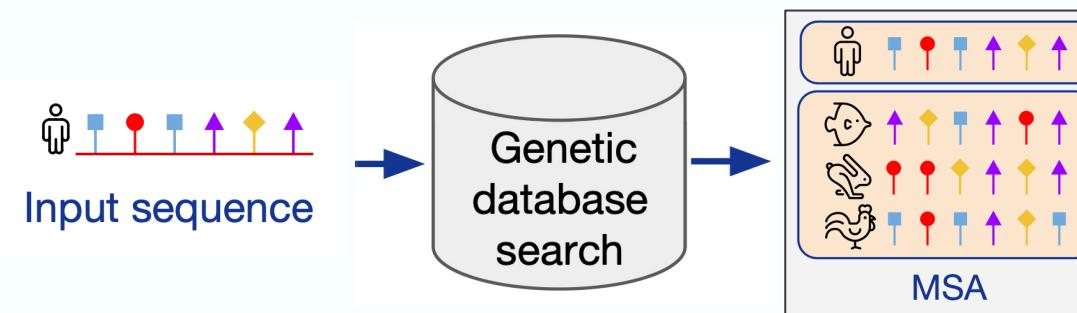
2024 NOBEL PRIZE IN CHEMISTRY

AlphaFold

How does it work?

Multiple sequence alignment

- MSAs are a way of associating amino-acid sequences of related proteins in different species.
- Provides information about *structure*:
 - Maintaining consistent structure is critical to protein function.
 - Amino acids critical to that structure are *conserved* across evolution.
 - Those that are not *vary* across evolution.
- AlphaFold exploits this: the input to the network is an MSA.

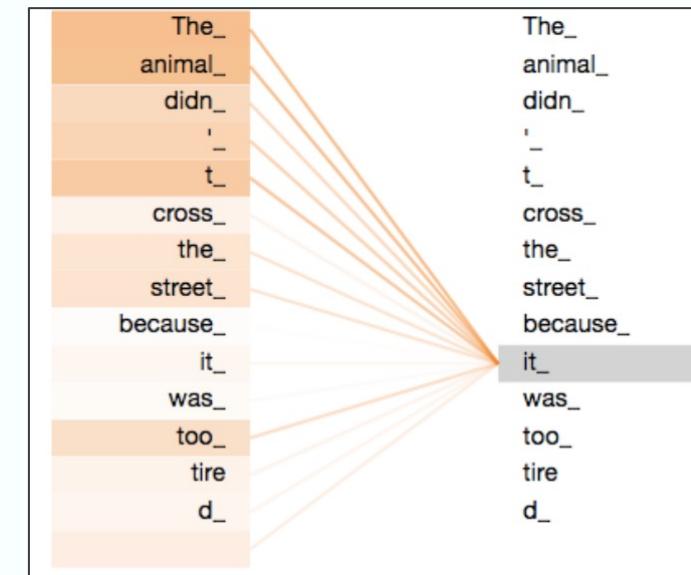


Evoformer

- The MSAs are processed by the *Evoformer* neural network module.
- A variant of the *transformer* neural network architecture (e.g., ChatGPT).



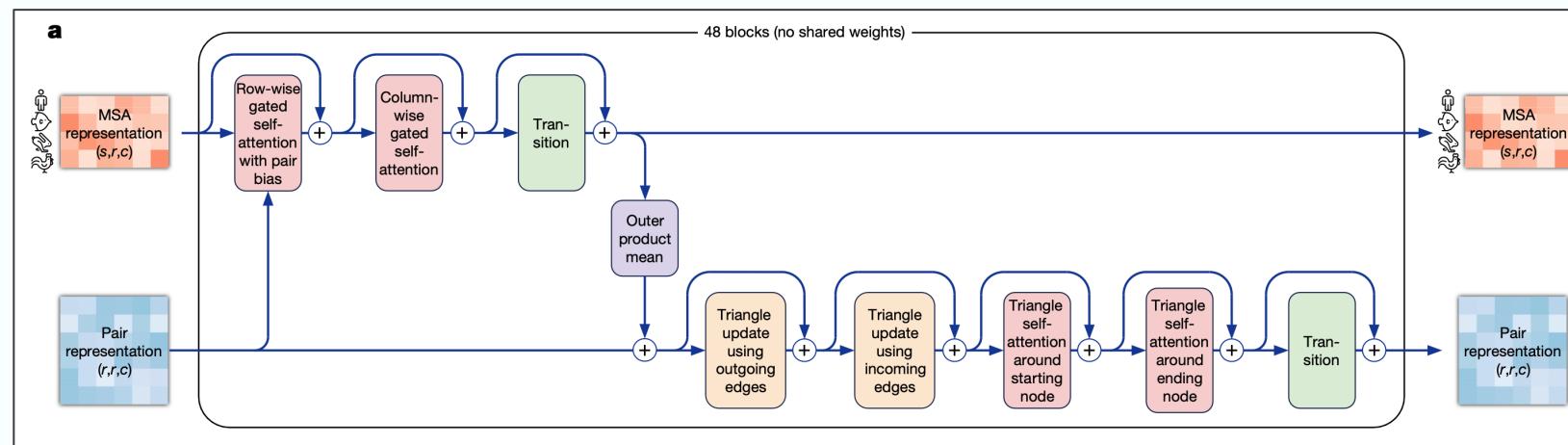
Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS'17). Curran Associates Inc., Red Hook, NY, USA, 6000–6010.



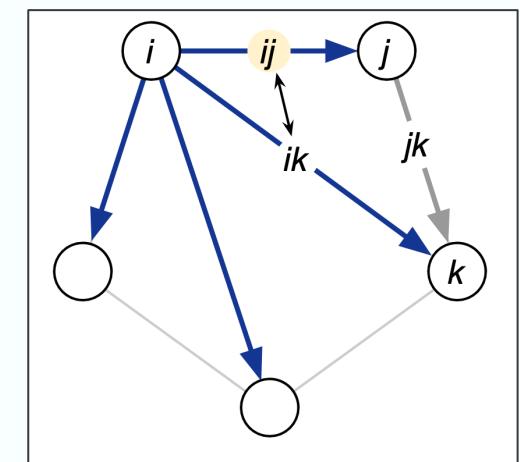
Self-attention in language generation

Evoformer

- Instead of attending to words, the network attends to *amino acids* in the MSA.
- The most important ingredient to AlphaFold's success.
- Specially-designed attention variants:
 - Triangle self-attention*—moving from pairwise communication ('A talks to B') to group dynamics ('A and B both talk to C, so what does that mean about A–C?')
 - Row-wise attention*—captures *within-sequence* context.
 - Column-wise attention*—captures *across-sequence* context.



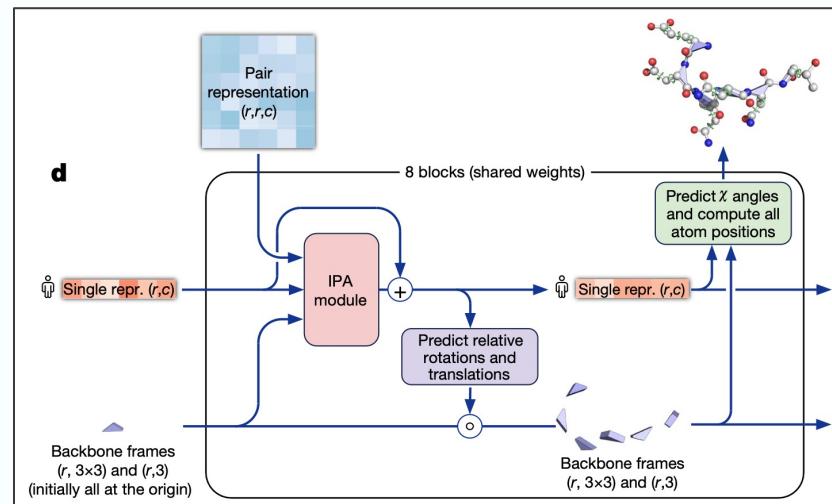
Evoformer architecture



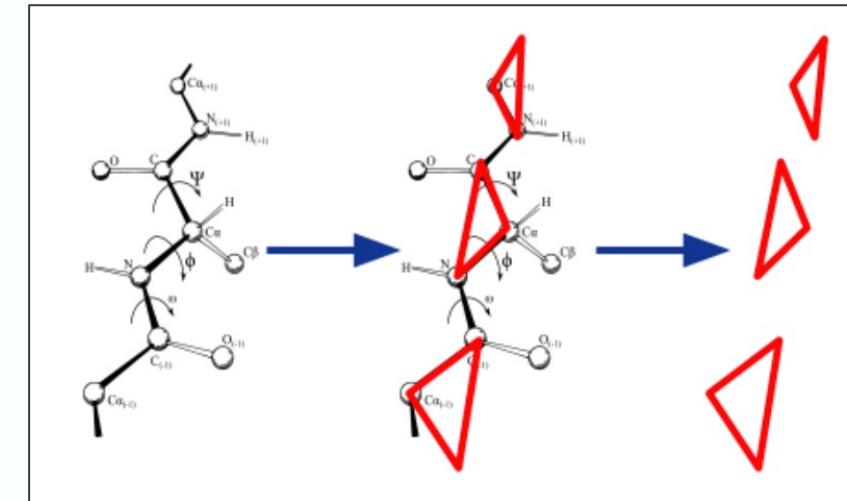
Triangle self-attention

Structure module

- The structure module takes the processed embeddings from the Evoformer and translates them into *3D geometry*.
- Each amino acid is represented by a *backbone-frame*:
 - Position in 3D space: (x, y, z) .
 - Orientation: (roll, pitch, yaw).



AlphaFold structure module



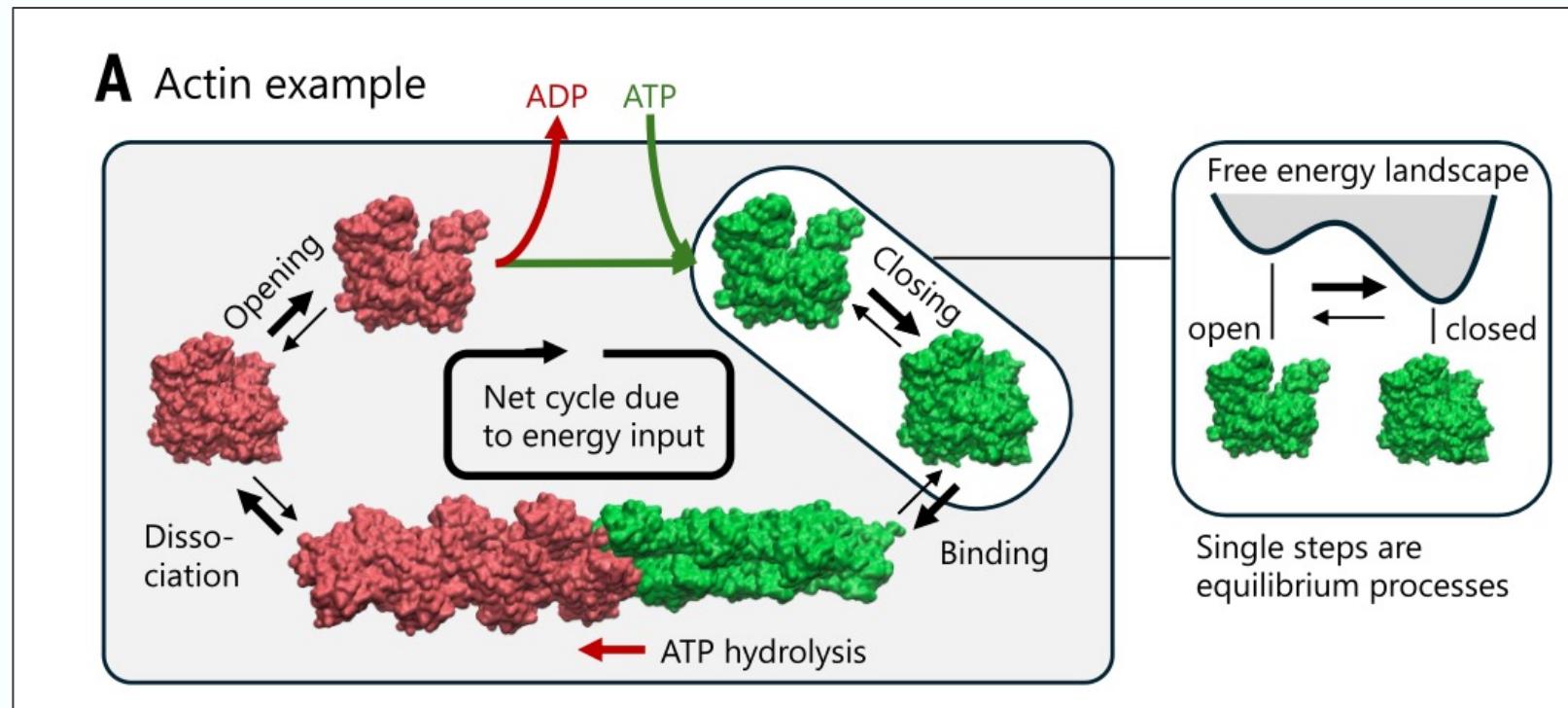
Backbone frame representation of amino acids

Structure is not enough

Equilibrium distributions

Protein equilibrium distributions

- Knowing the folded state of a protein is not everything.
- Protein *function* is often mediated by *transitions between multiple conformations*.



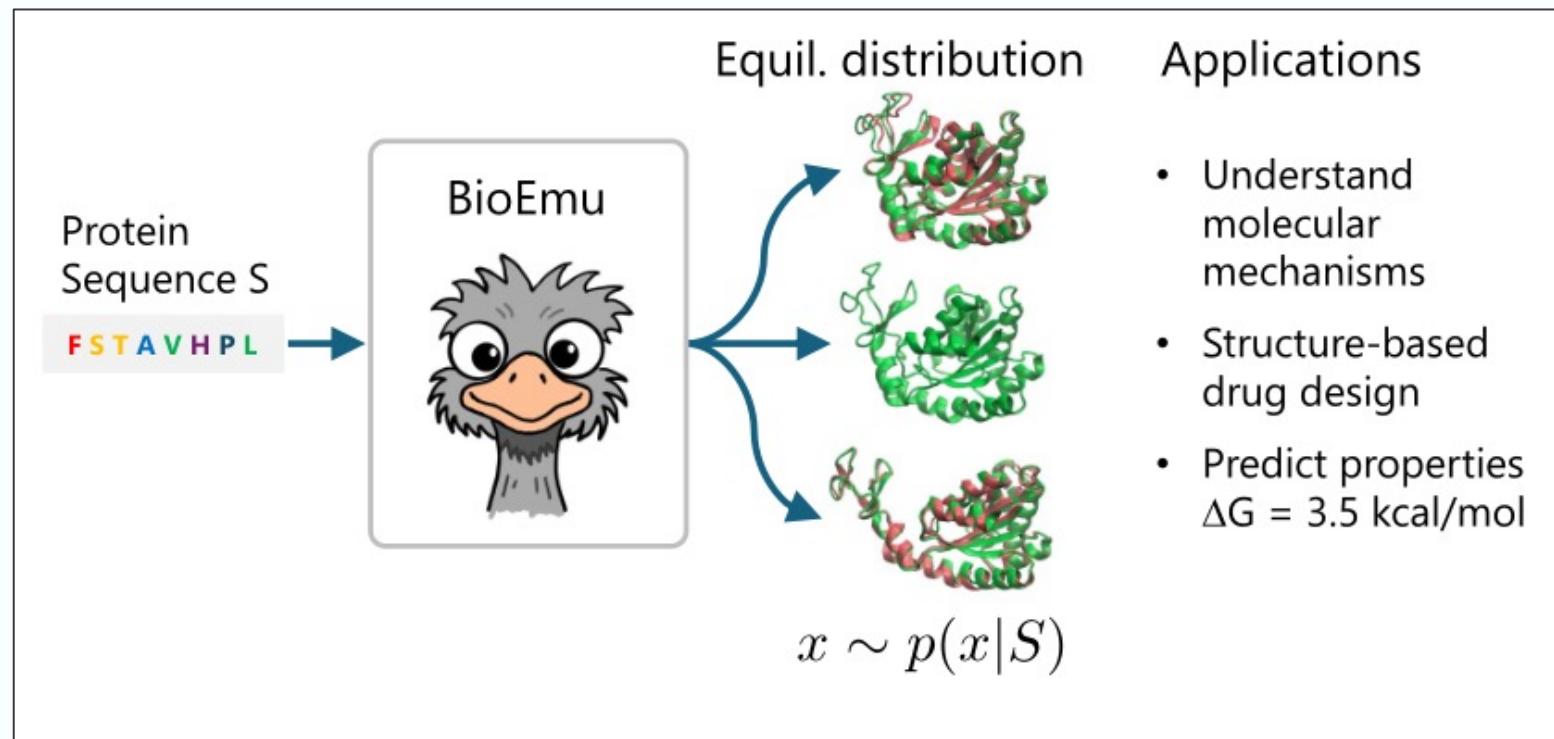
Conformational change drives function

Protein equilibrium distributions

- BOLTZMANN DISTRIBUTION: $p(x) = \exp(-U(x)/k_B T) / Z$.
 - x is the 3D conformation of the protein.
 - $p(x)$ is its probability density.
 - $U(x)$ is the potential energy of a conformation.
 - k_B is Boltzmann's constant.
 - T is the temperature.
 - Z is a normalizing constant.
- *If we can sample from the Boltzmann distribution, we can describe every possible state of a protein and its relative probability.*

Going beyond AlphaFold

- BioEmu goal: approximate sampling from Boltzmann distribution *without* expensive molecular dynamics simulations.



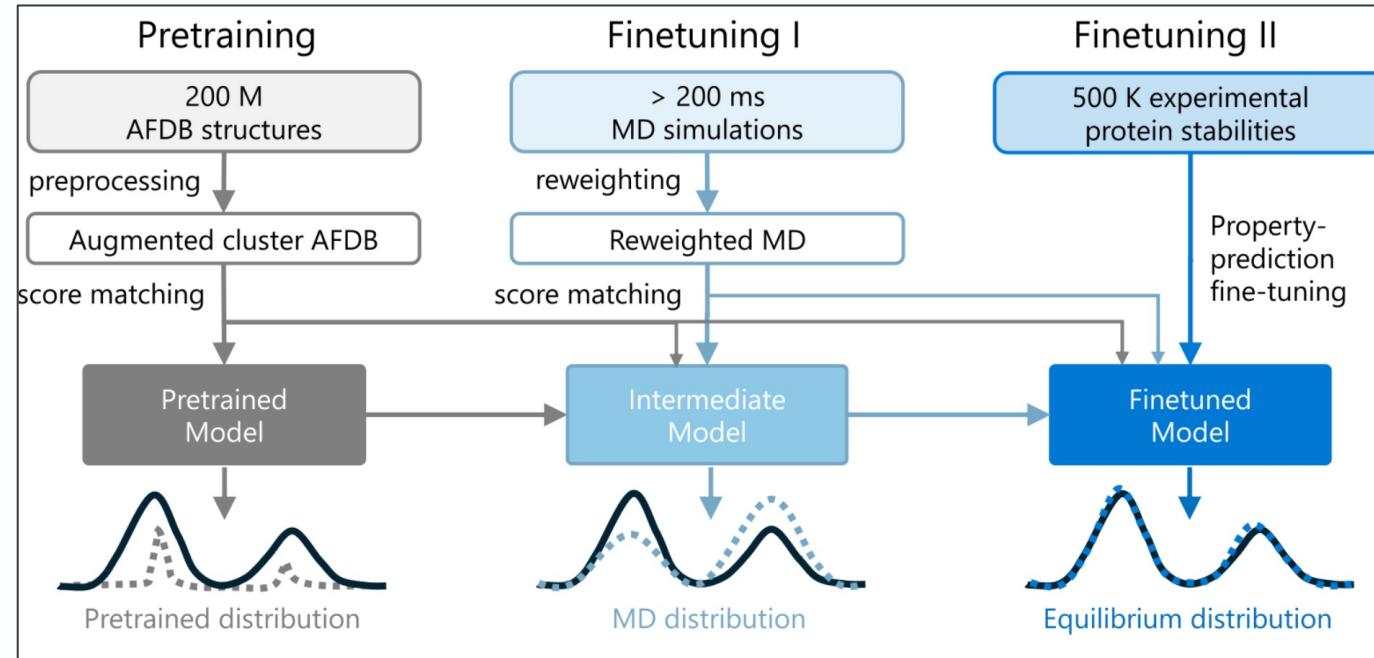
BioEmu in a nutshell

BioEmu

Training data

Sources of data

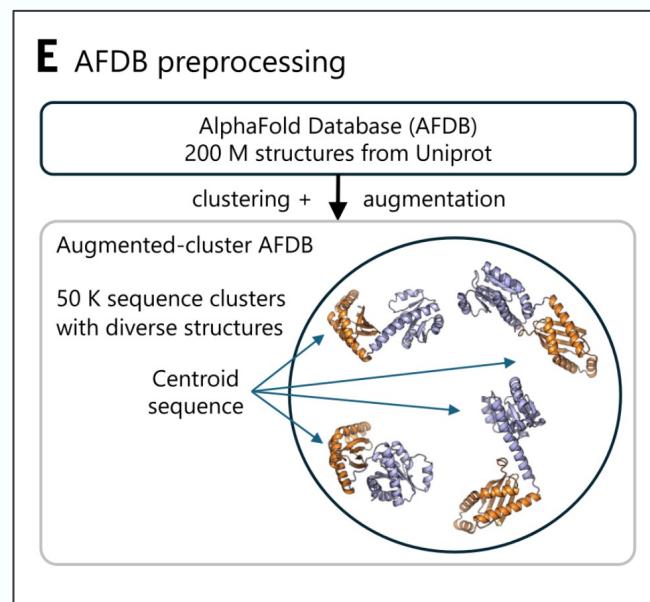
- First challenge: *data for equilibrium distributions is hard to come by.*
- There is no “PDB” for distributional data.
- Our approach:



Three-stage approach for BioEmu training data

Sources of data

- Pretraining on diversified AlphaFold *predictions*.
 - Download 200 million protein structure predictions from AlphaFold Database.
 - Cluster by sequence similarity.
 - Group *diverse structures with similar sequences* together.
- CONSEQUENCE: BioEmu is forced to learn *conformational variability*.



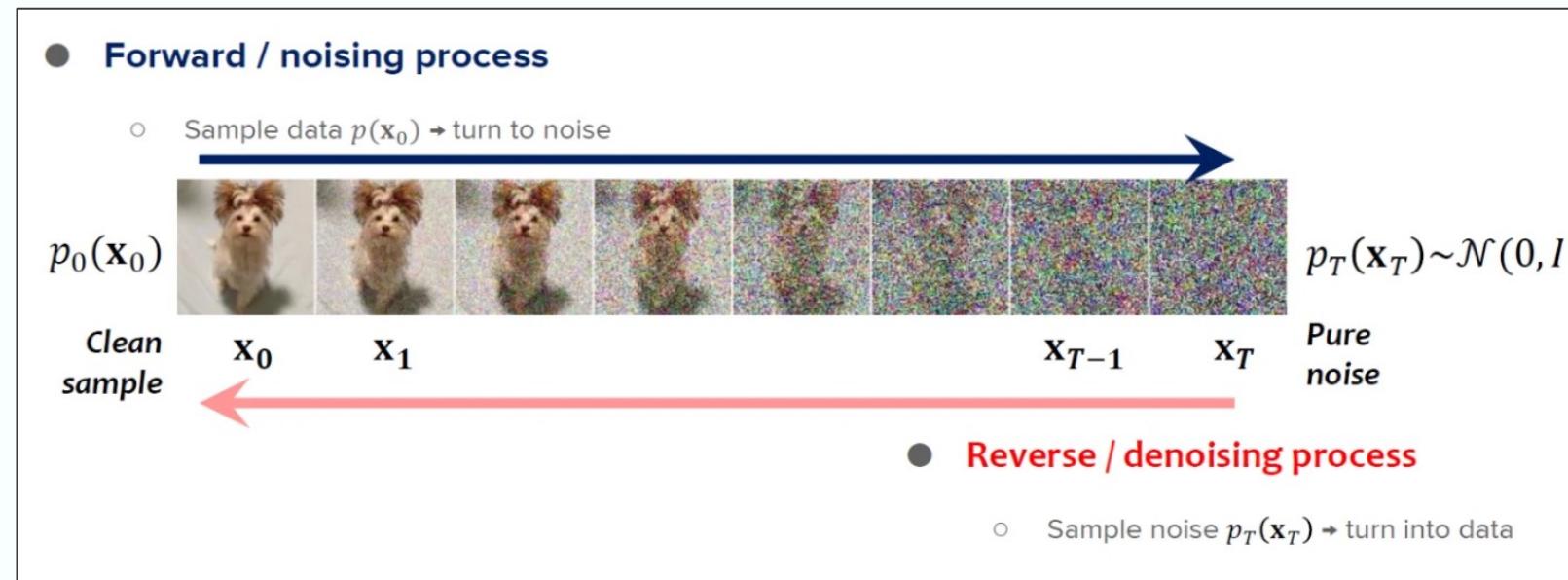
Pretraining with
diversified AFDB
structure clusters

BioEmu

Diffusion model

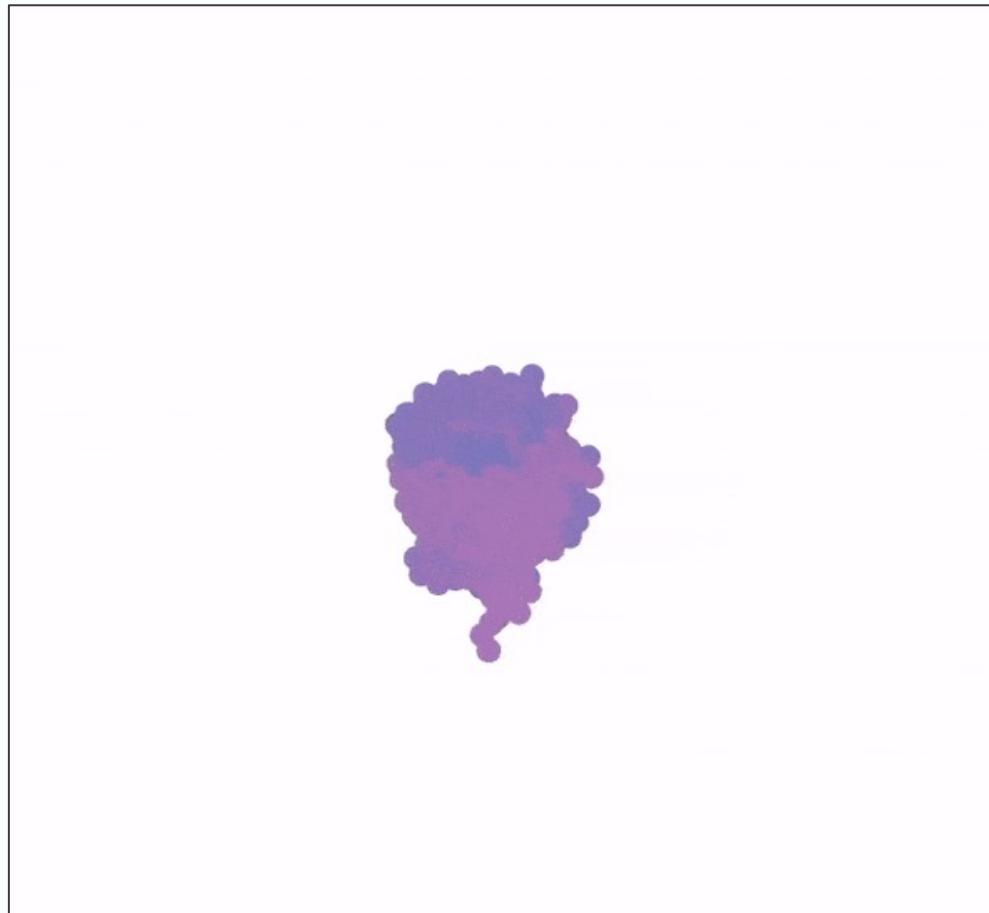
Diffusion models

- AlphaFold has *one output structure for one input sequence.*
- DIFFUSION MODELS output *distributions over data:*
 - Random noise input is converted to structures.
 - Variability in the noise translates to variability in the structures.
- Ideal candidate for a Boltzmann distribution emulator.



Diffusion models learn how to map random noise to data

Diffusion models

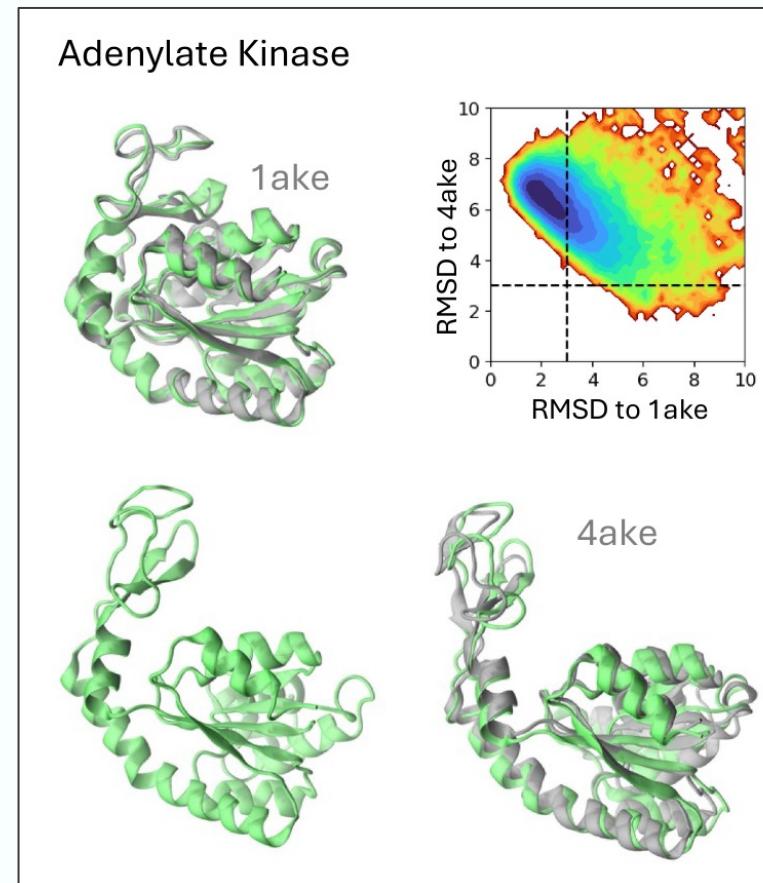


EXAMPLE: diffusion model generating protein structure

BioEmu

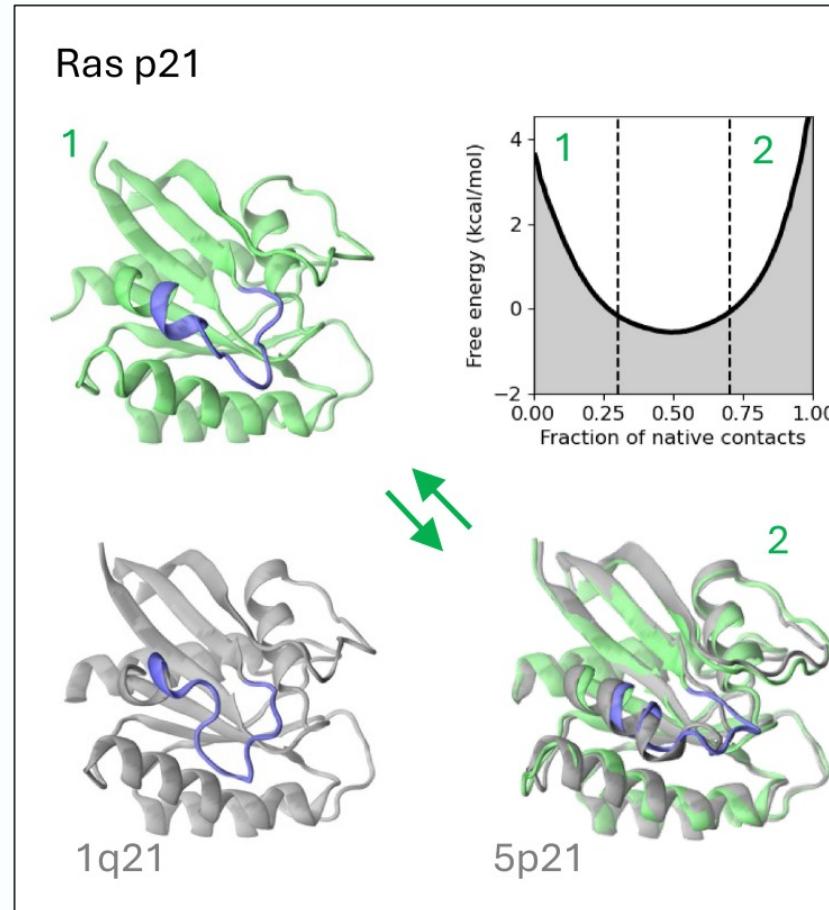
Capabilities on unseen proteins

Domain motions



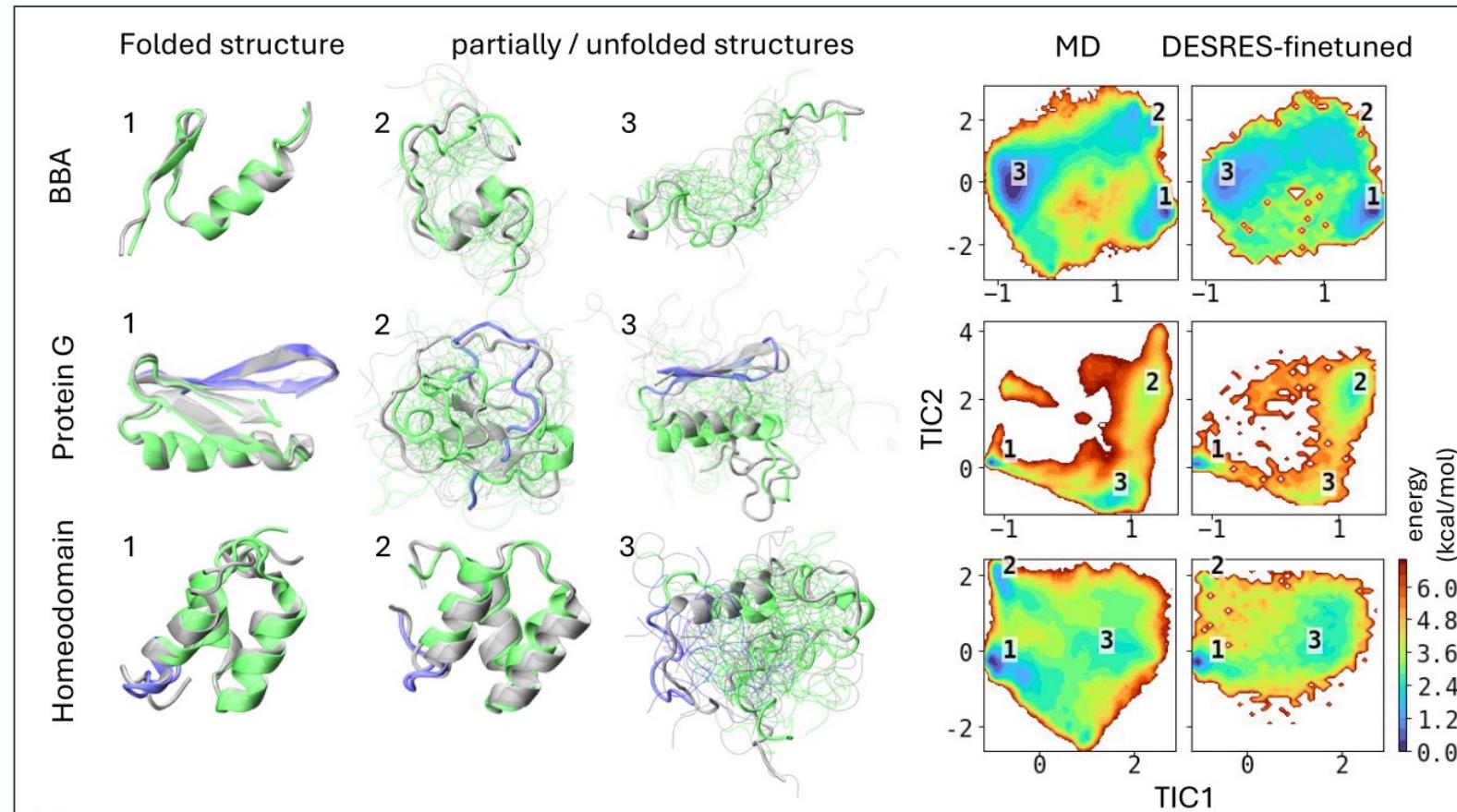
BioEmu predicts large-scale domain motions

Local unfolding



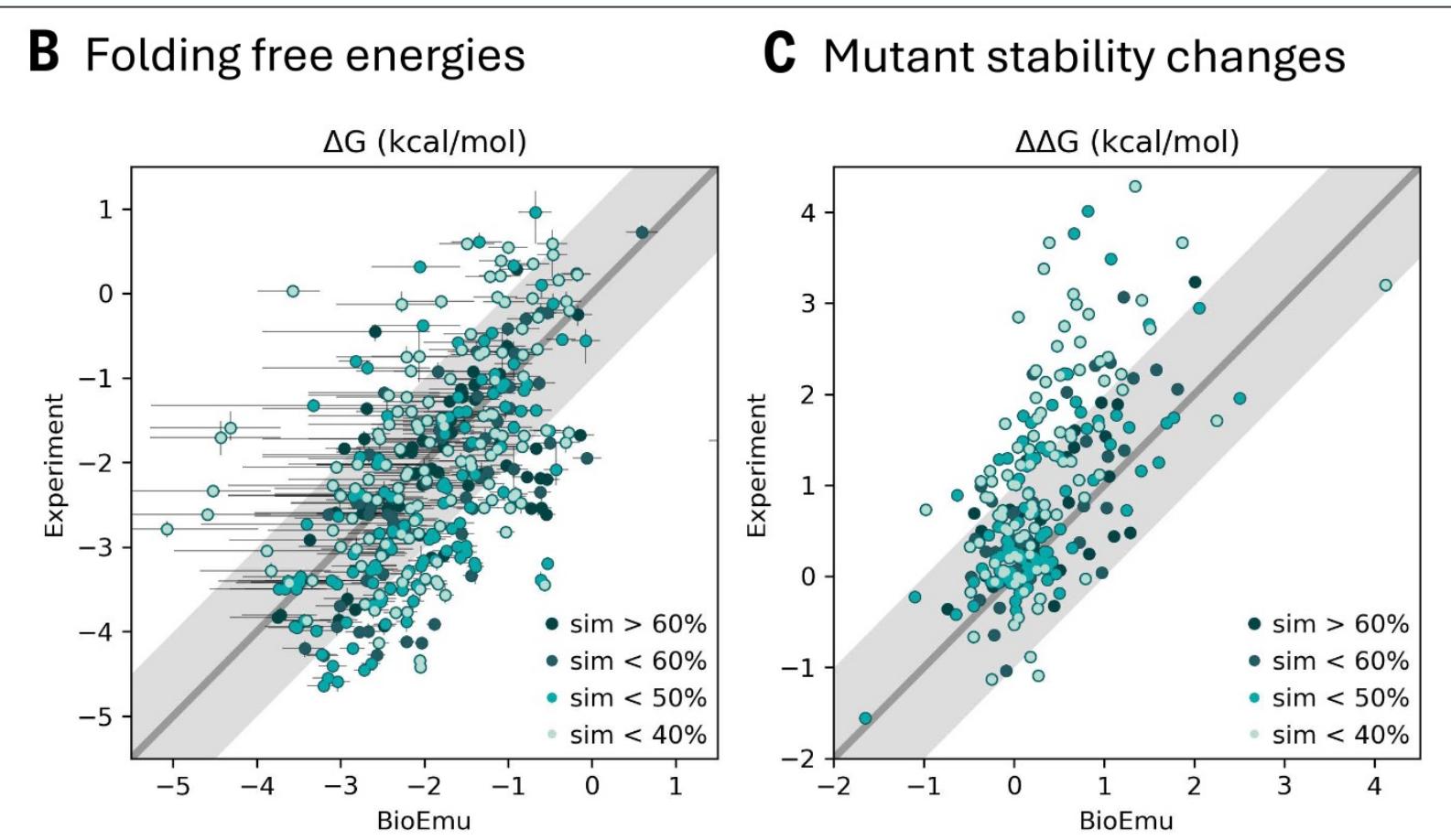
BioEmu predicts local unfolding transitions (*alpha helix destabilizes*)

Comparison with molecular dynamics



Sampling from BioEmu yields similar distributions of structures as lengthy MD simulations

Free-energy predictions



BioEmu samples reflect relative stability of protein mutants

Using BioEmu on ColabFold

Try BioEmu for free today:

- Runs in browser; no software installation required.
- Just input a sequence and download the sampled structures.

Biomolecular Emulator (BioEmu) in ColabFold

[BioEmu](#) is a framework for emulating biomolecular dynamics and integrating structural prediction tools to accelerate research in structural biology and protein engineering. This notebook uses BioEmu with ColabFold to generate the MSA and identify cluster conformations using Foldseek.

For more details, please read the [BioEmu Preprint](#).



To run

Either run each cell sequentially, or click on `Runtime -> Run All` after choosing the desired sampling config

Sample with following config

```
▶ • sequence : Monomer sequence to sample  
sequence: " MADQLTEEQIAEKFKEAFSLFDKDGDTITTKELGTVMRSLGQNPTEAEIQLQDMINEVDADGNGTIDFPEFLTMMARKMKDTSSEEIREAFRVFDKGNGYISAAELF "  
• num_samples : Number of samples requested  
num_samples: 10
```

<https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/BioEmu.ipynb>

Thanks for listening!

Questions?