# Generative and Geometric Deep Learning in Computational Biology

Darin Tsui

# Computational biology

Intersection of computer science, biology, and big data →

Applied mathematics

Genetics

Chemistry

# Why machine learning?

Images:
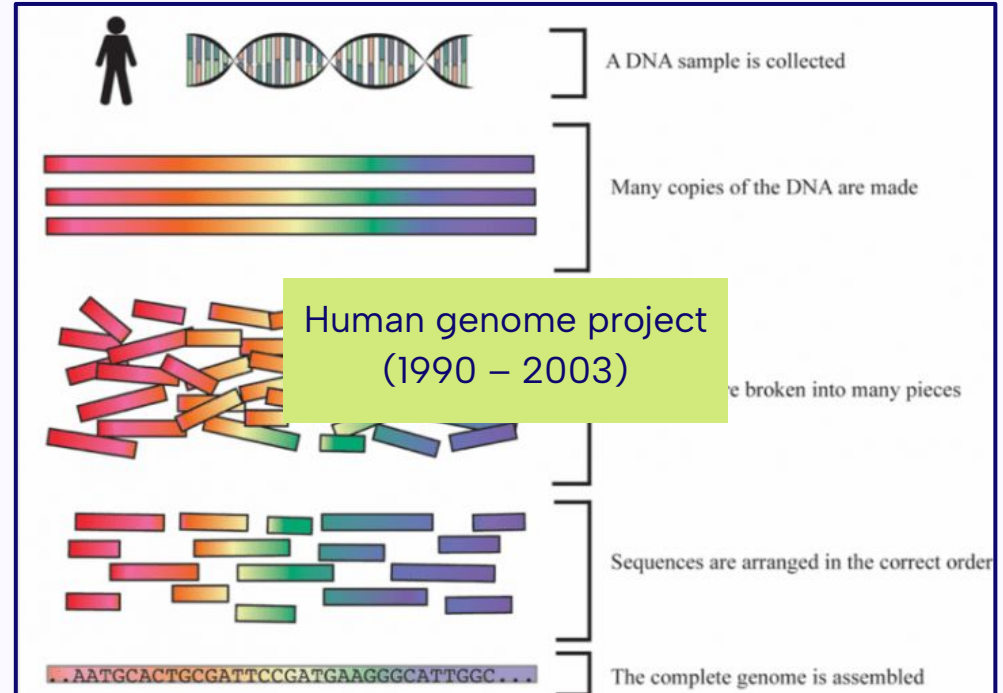**high-dimensional** data



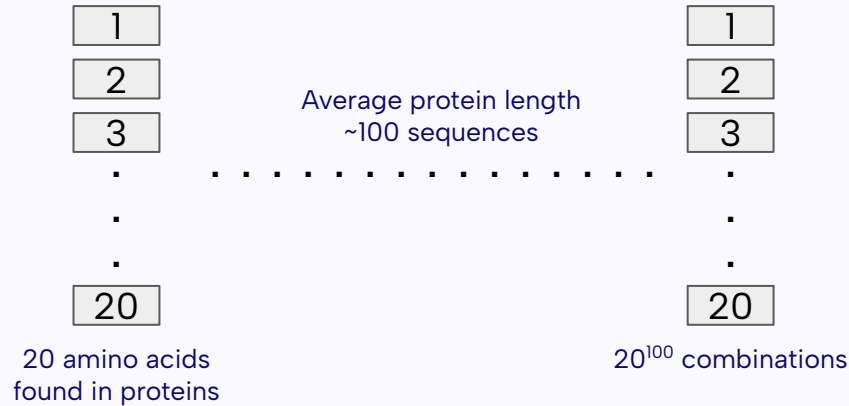MNIST (1994)

28 x 28 = **784 features**

# Why machine learning?

3,200,000,000 nucleotides

100,000 different proteins

20,000 – 25,000 genes



A DNA sample is collected

Many copies of the DNA are made

Human genome project
(1990 – 2003)

...re broken into many pieces

Sequences are arranged in the correct order

. . .AATGCACTGCGATTCCGATGAAGGGCATTGGC. . .  The complete genome is assembled

# How many features?



1
2
3
.
.
.
20

Average protein length
~100 sequences

1
2
3
.
.
.
20

20 amino acids
found in proteins

$20^{100}$ combinations

Over 98% of relevant protein variants still have unknown consequences*

Frazer, J. *et al*. *Nature* (2021)

# Why generative AI?

Labels are:

–   Biased

–   Variable quality

We can use generative AI to make predictions on unknown sequences using unlabelled data



A DNA sample is collected

(1990 – 2003)

...re broken into many pieces

Sequences are arranged in the correct order

. .AATGCACTGCGATTCCGATGAAGGGCATTGGC. . .

The complete genome is assembled

# Autoregressive models for protein design

# Autoregressive models for protein design

$$p\left(\mathbf{x}|\boldsymbol{\theta}\right) = p(x_1|\boldsymbol{\theta})\prod_{i=2}^{L} p(x_i|x_1,\ldots,x_{i-1};\boldsymbol{\theta})$$

**x** = sequence

**θ** = constraints for functional sequences

p(**x**|**θ**) = probability of sequence generation given evolution

Training set to predict masked values of amino acids

Shin, J. *et al. Nature Comm* (2021)

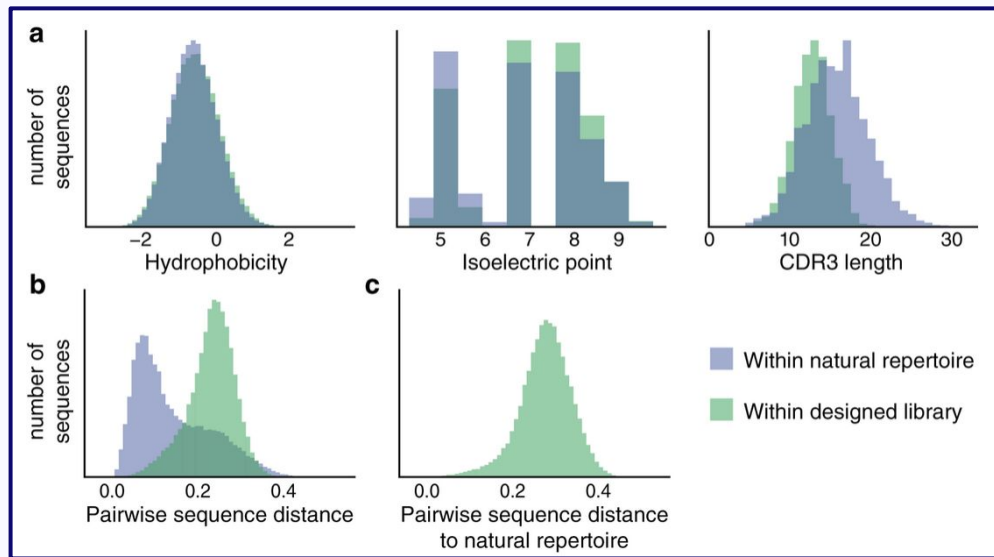# Autoregressive models for protein design

$$\log \frac{p(\mathbf{x}^{\mathrm{Mutant}} | \boldsymbol{\theta})}{p(\mathbf{x}^{\mathrm{Wild-type}} | \boldsymbol{\theta})}$$

Predicting mutation effects using the
**log–ratio of likelihoods**
- Summed cross–entropy between true vs predicted amino acids at each position, conditioned on previous amino acids

Shin, J. *et al. Nature Comm* (2021)

# Autoregressive models for protein design



Shin, J. *et al. Nature Comm* (2021)

# Autoregressive models for protein design

Takeaway: amino acids behave like **text**

Autoregressive likelihood = context–dependent prediction

Do not rely on "word alignment", unlike other models

Potential **cons**:

Reliant on massive amounts of data (~1.2 million sequences)

Bottlenecked by sequences/ conditional probabilities

# Alignment-based methods

## Autoregressive Likelihood

Input sequences are of variable length

Raw sequences don't capture evolutionary information

MTAIIKEIVSRNKRRYQED

## Multiple Sequence Alignment (MSA)

Fixed length

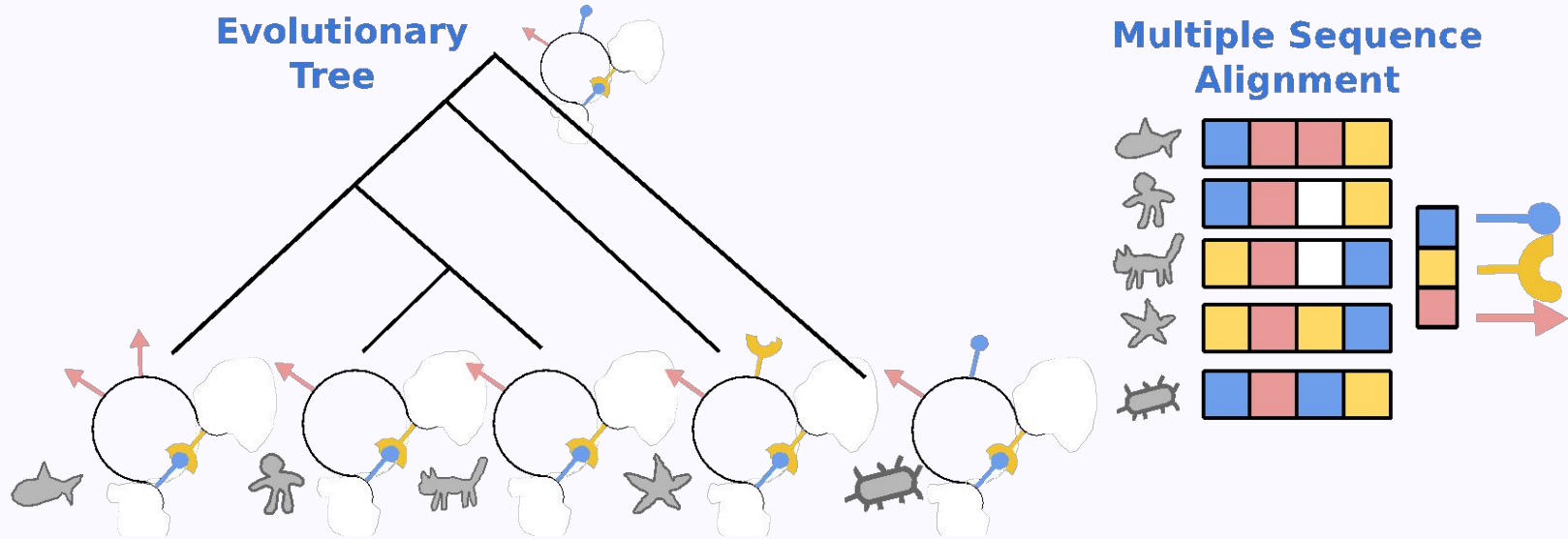Sequences are clustered in "**protein families**"

MTAIIKEIVSRNKRRYQED
MTAIIKEIVTTNKRRTQED
MTBIIKEIVSCNKRRTQED

# Multiple Sequence Alignment (MSA)



Rao, R. (2021)

# Multiple Sequence Alignment (MSA)

**Evolutionary Tree**

**Multiple Sequence Alignment**

Query protein database for related sequences

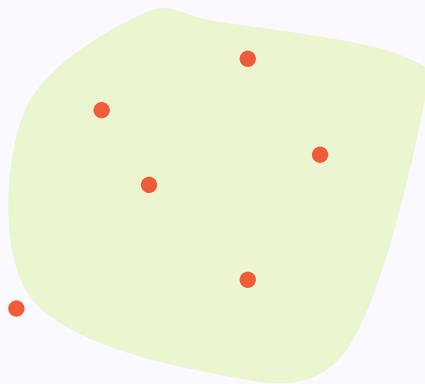Align with heuristic based on edit distance

Rao, R. (2021)

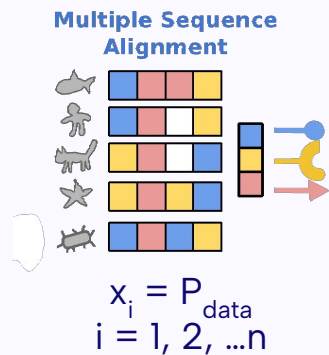# MSA information to learn data distributions
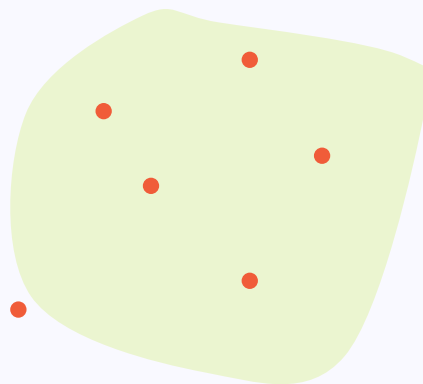
$x_i = P_{data}$
$i = 1, 2, ...n$

All images

Model family $\theta \subseteq M$

# MSA information to learn data distributions

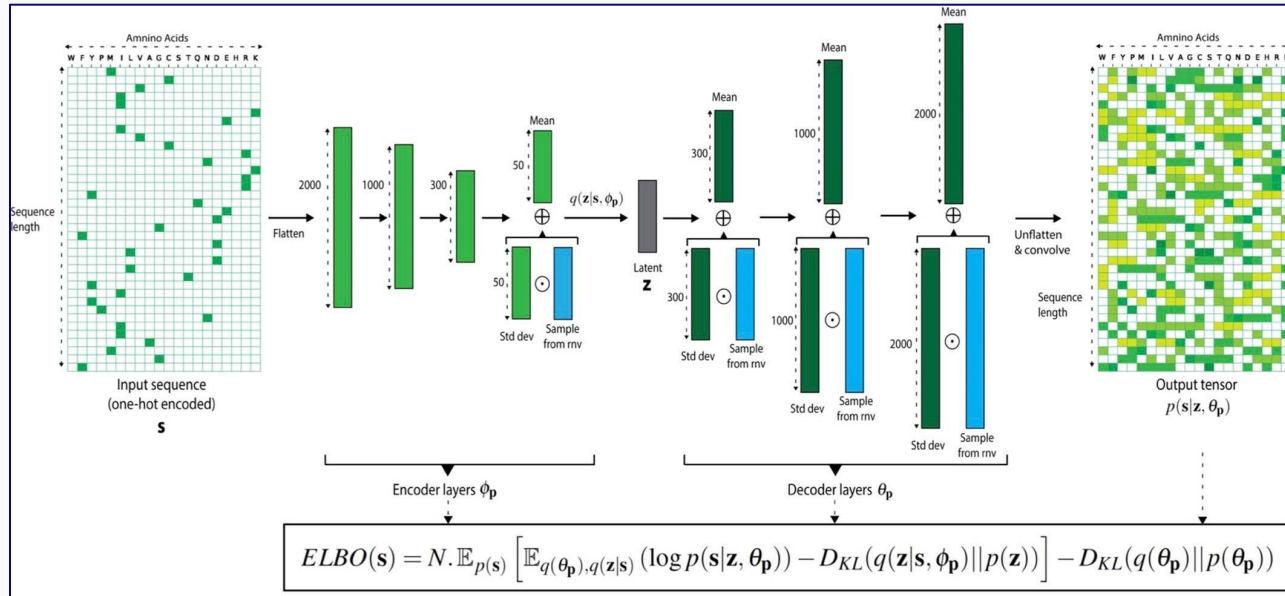**Multiple Sequence Alignment**

$x_i = P_{data}$
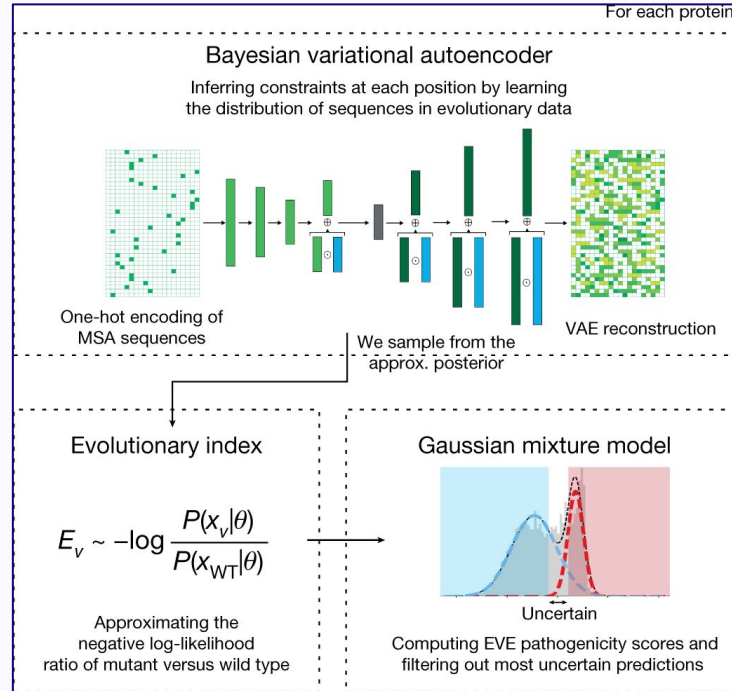$i = 1, 2, ...n$

$20^n$ sequences

Model family $\theta \subseteq M$

Rao, R. (2021)

# Variational Autoencoders (VAEs) for variant prediction



VAE learns a **protein family** distribution from MSA data

Frazer, J. *et al*. *Nature* (2021)

# Variational Autoencoders (VAEs) for variant prediction



Frazer, J. *et al*. *Nature* (2021)

# Variational Autoencoders (VAEs) for variant prediction



Evolutionary Index (EI) = ELBO(**w**) - ELBO(**s**)

- **w** = wild-type (protein family "parent")
- **s** = mutated sample

Frazer, J. *et al. Nature* (2021)

# Variational Autoencoders (VAEs) for variant prediction

Takeaway: MSA information can be learned as a **distribution**

MSA information generally implies better generalization

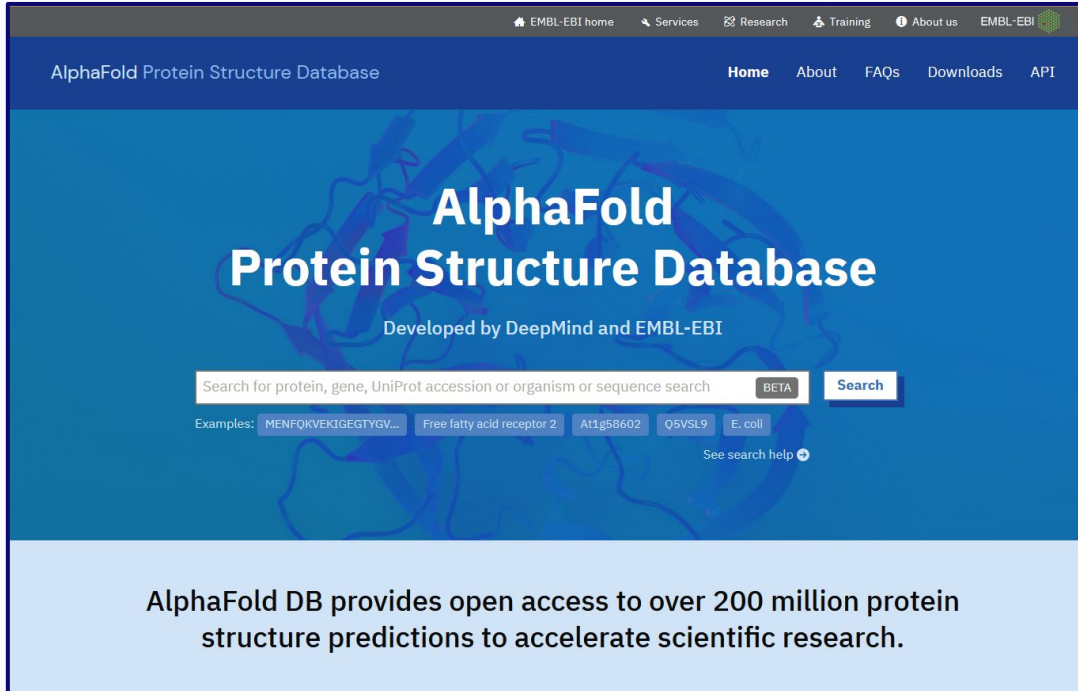**Significantly** less input data than in autoregressive models

Potential **cons**:

Unclear what the VAE is exactly learning (how to validate?)

How *precise* is the VAE with the limited data its given?

Frazer, J. *et al. Nature* (2021)

# Geometric deep learning for protein structure prediction

# Geometric deep learning for protein structure prediction

# Why predict protein structures?



**Extremely** hard to determine experimentally. Only roughly about 200k protein structures available at the Protein Data Bank (PDB)

Bates, R. *et al. DeepMind* (2021)

# Intuition meets machine learning



Proteins conform to physics to go from a 2D sequence to a 3D structure. Evolutionary history gives us information about structure.

Bates, R. *et al. DeepMind* (2021)

# Intuition meets machine learning



Convolutional Networks (e.g. computer vision)
- data in regular grid
- information flow to local neighbours
- AlphaFold 1

Recurrent Networks (e.g. language)
- data in ordered sequence
- information flow sequentially

Graph Networks (e.g. recommender systems or molecules)
- data in fixed graph structure
- information flow along fixed edges

Attention Module (e.g. language)
- data in unordered set
- information flow dynamically controlled by the network (via keys and queries)

**Graph networks** and **attention modules** best reflect our understanding of physics and geometry

Bates, R. *et al. DeepMind* (2021)

# Intuition meets machine learning

Input: Amino acid sequence

Output: 3D structure

Model should reflect our understanding of physics and geometry

Amino acid positions are de-emphasized

Amino acids close in 3D position should communicate

Network learns a graph

Bates, R. *et al. DeepMind* (2021)

# AlphaFold Architecture

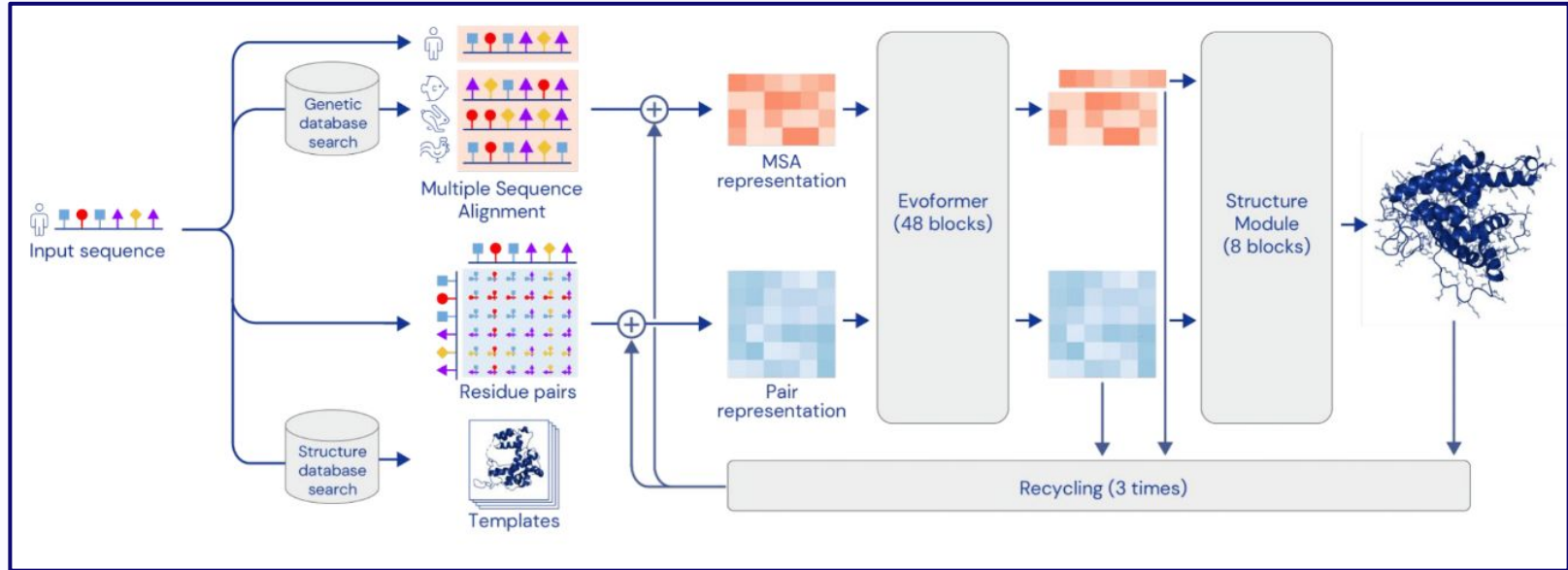

Bates, R. *et al. DeepMind* (2021)

# Evoformer Architecture



Jumper, J. *et al*. *Nature* (2021)

# Evoformer - Triangular Attention



Purpose of pair representation is to encode distances between amino acids (i,j)

- Constrain these pairwise representations by the triangle inequality
- Allows for updating that's more consistent with 3D structure

Jumper, J. *et al*. *Nature* (2021)

# Structure Module: Invariance and Equivariance



End-to-end folding instead of gradient descent

Protein backbone = gas of 3-D rigid bodies (chain is learned!)

Invariant Point Attention (IPA): update seq representation without affecting 3D positions. Augments local frames of reference without affecting the global frame of reference

Equivariant update using updated sequence representation

Jumper, J. *et al*. *Nature* (2021)

# A collection of innovations



a

Test set of CASP14 domains | Test set of PDB chains
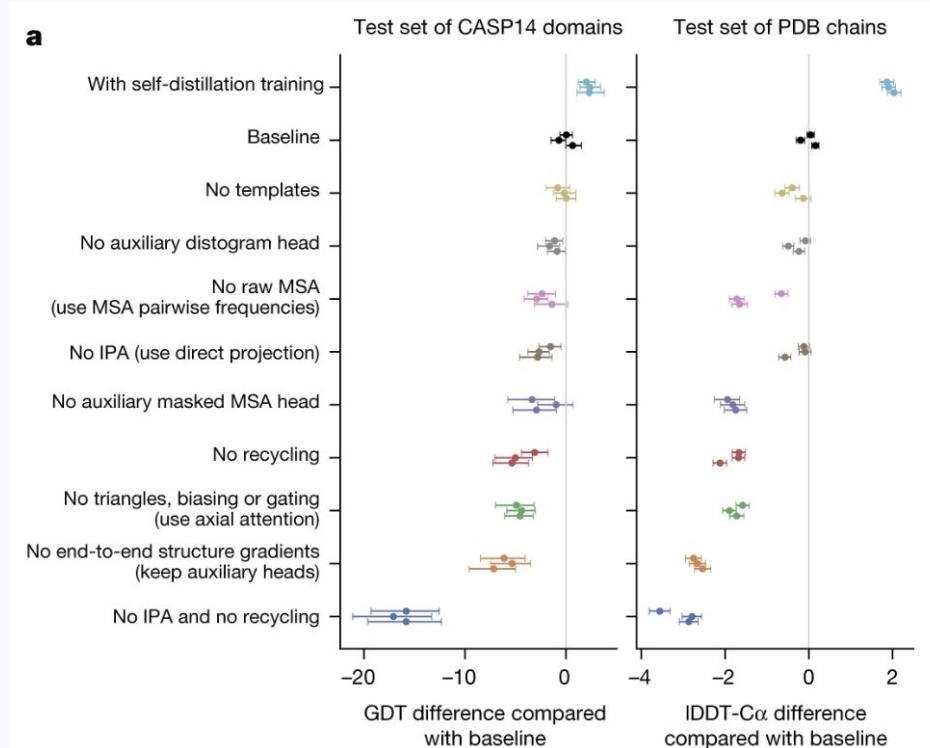
With self-distillation training
Baseline
No templates
No auxiliary distogram head
No raw MSA (use MSA pairwise frequencies)
No IPA (use direct projection)
No auxiliary masked MSA head
No recycling
No triangles, biasing or gating (use axial attention)
No end-to-end structure gradients (keep auxiliary heads)
No IPA and no recycling

GDT difference compared with baseline

IDDT-Cα difference compared with baseline

Jumper, J. *et al.* *Nature* (2021)

# AlphaFold for Structure Prediction

Takeaway: no one change made structure prediction possible

Emphasis on geometric constraints

A need for more "smarter" neural networks

Potential **cons**:

Constrained by high-depth MSA data

Slow inference time

Jumper, J. *et al*. *Nature* (2021)

# What's the future of computational biology?

Need for **interpretable** algorithms

Need for **scalable** algorithms

Need for **smarter** algorithms

# Thank you!