# Protein Structure Prediction

BCH394P/364C Systems Biology/Bioinformatics
April 1, 2025
Daryl Barth

# Today's Goals

- Motivation and bit of history on protein structure prediction
- CASP
- AlphaFold
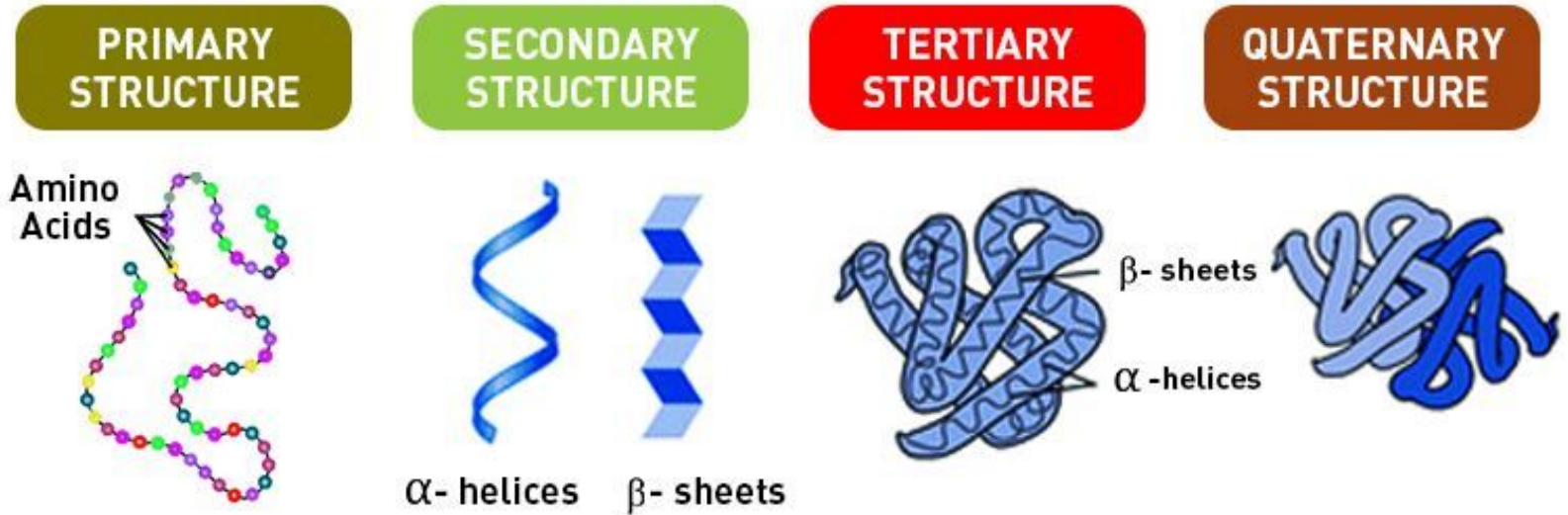- Metrics you need to know
- ColabFold
- Demo!

# Why?

# The four levels of protein structure

# Is an amino acid sequence all you need?

- Anfinsen's dogma:
- The Protein Folding Problem
  - What is the folding code?
  - What is the folding mechanism?
  - Can we predict a native protein structure from its primary, amino acid sequence?

A protein's native structure stands for a free energy minimum determined by its amino acid sequence...
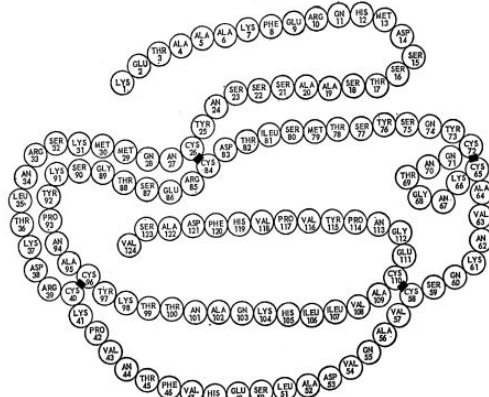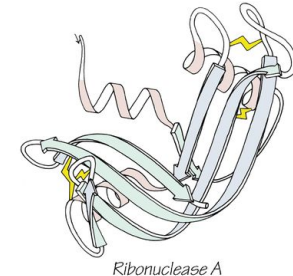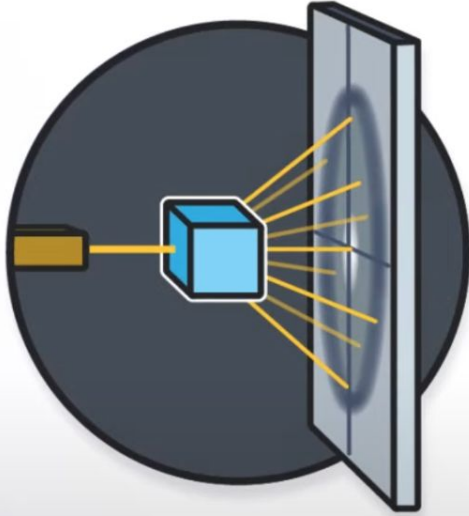


Fig. 1. The amino acid sequence of bovine pancreatic ribonuclease (50).





Ribonuclease A

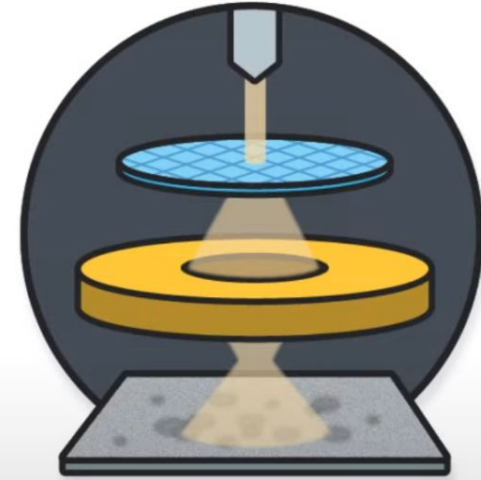# Experimental Methods for Determining Protein Structure



X-Ray crystallography

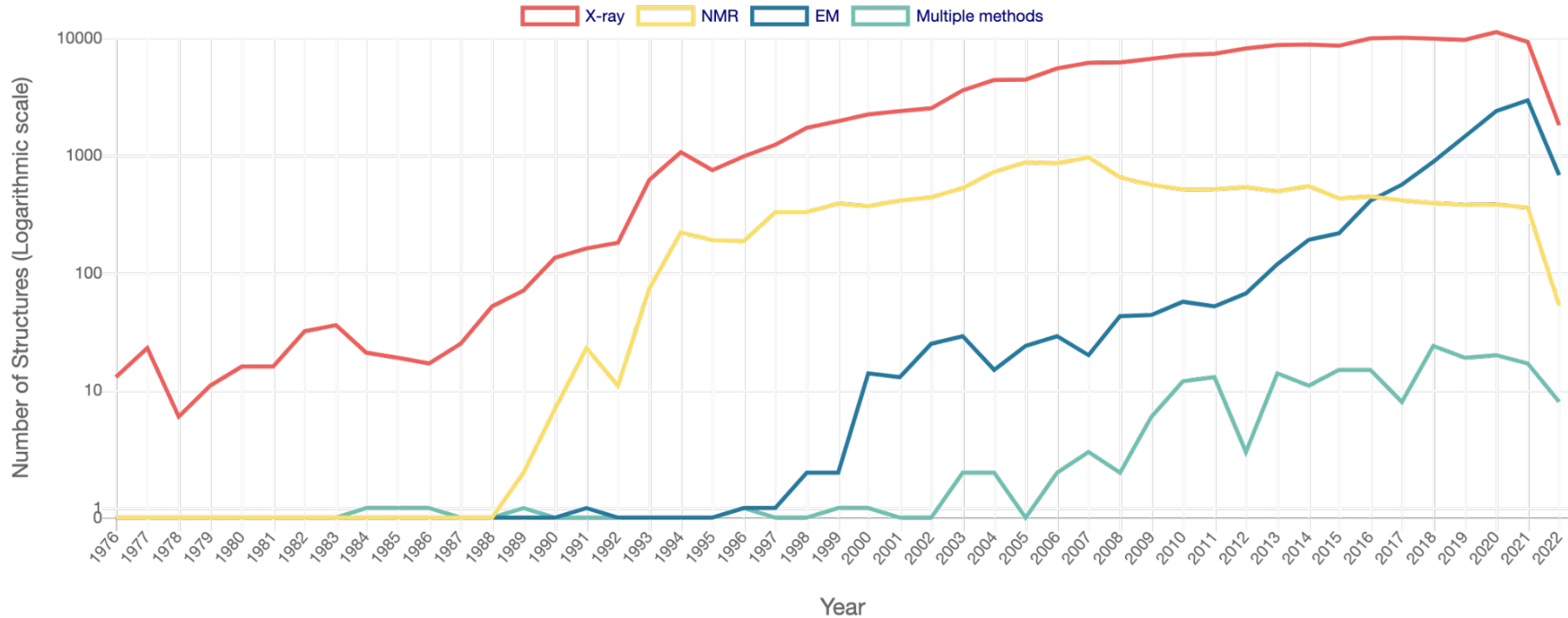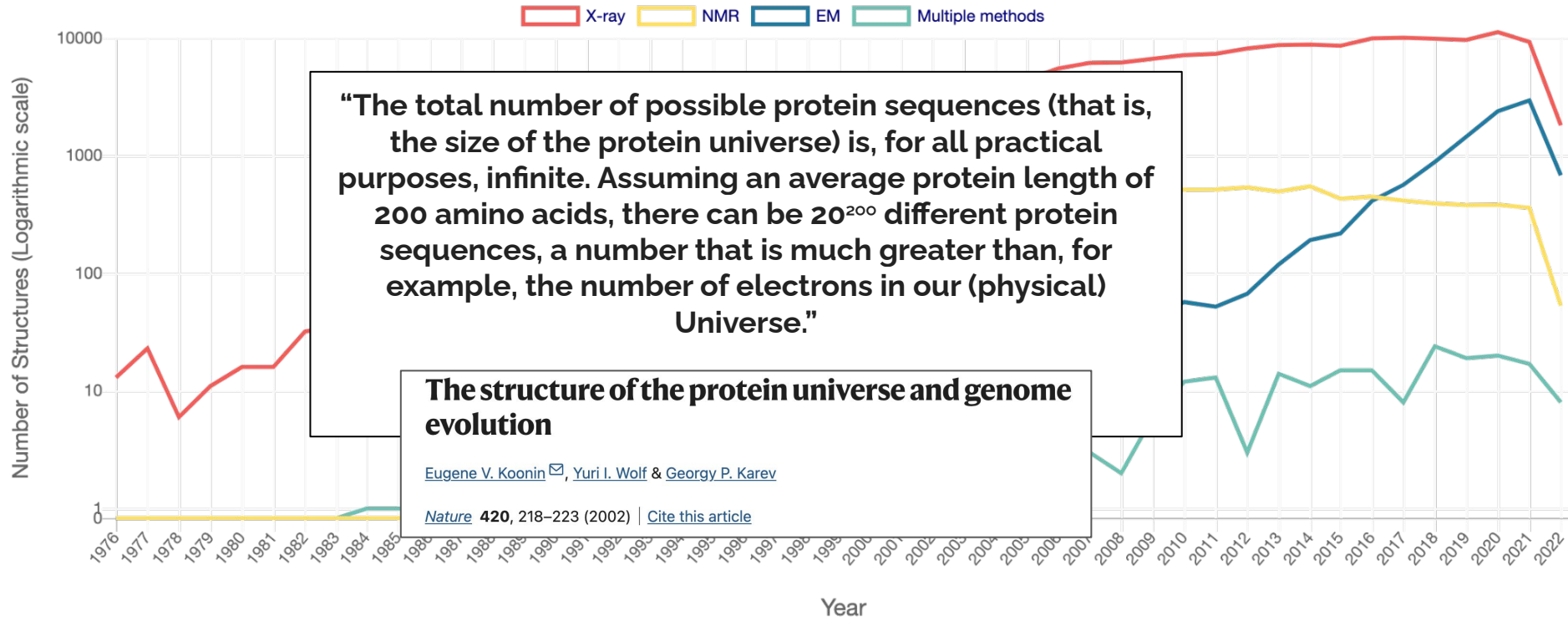Nuclear magnetic resonance spectroscopy

Cryoelectron microscopy

# How many structures have been determined?

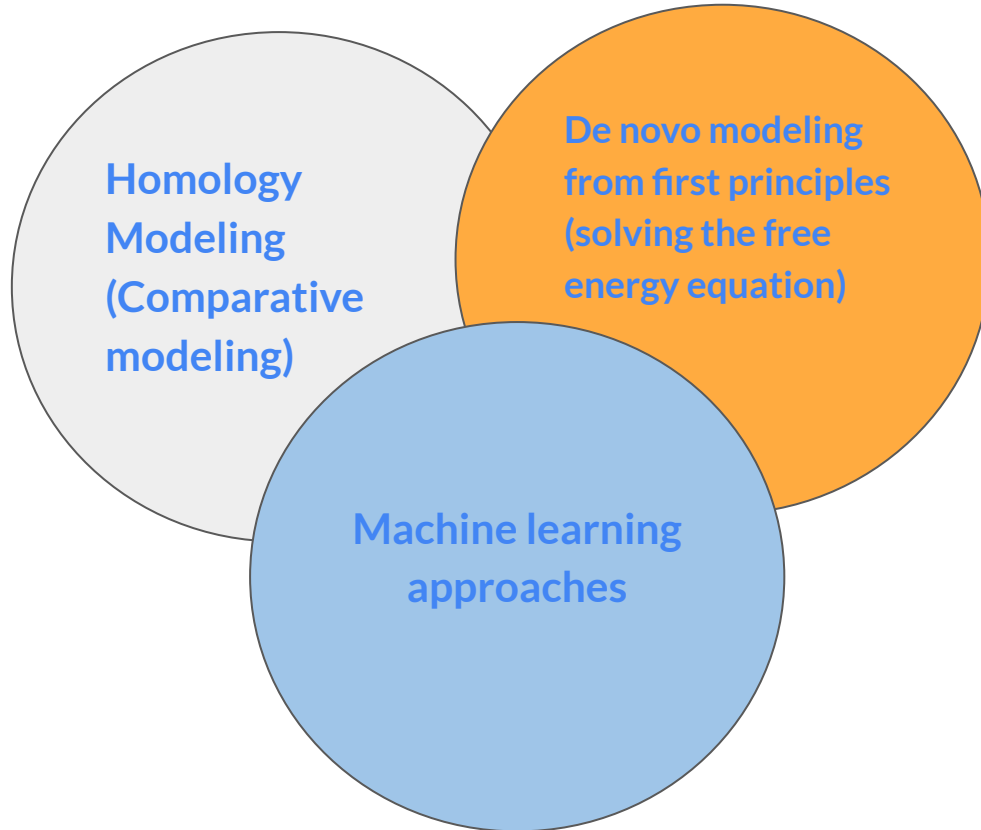# How many structures have been determined?



"The total number of possible protein sequences (that is, the size of the protein universe) is, for all practical purposes, infinite. Assuming an average protein length of 200 amino acids, there can be $20^{200}$ different protein sequences, a number that is much greater than, for example, the number of electrons in our (physical) Universe."

**The structure of the protein universe and genome evolution**

Eugene V. Koonin ✉, Yuri I. Wolf & Georgy P. Karev

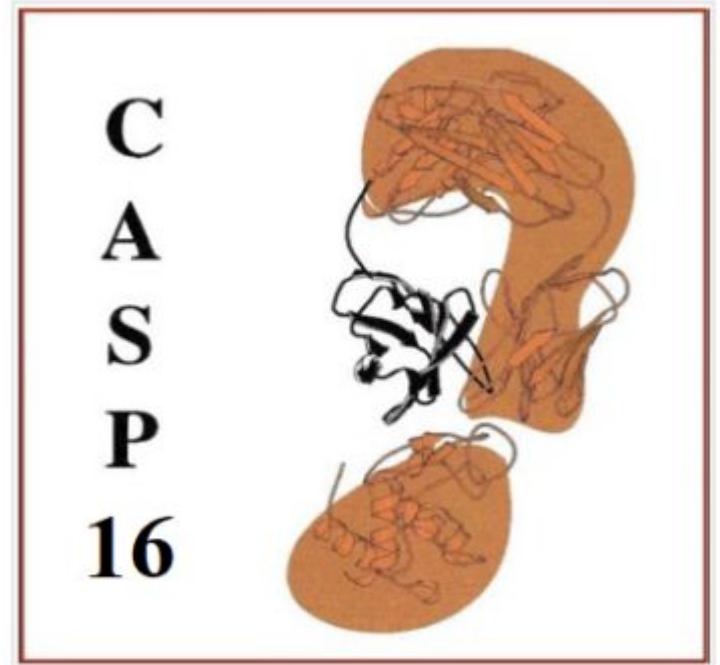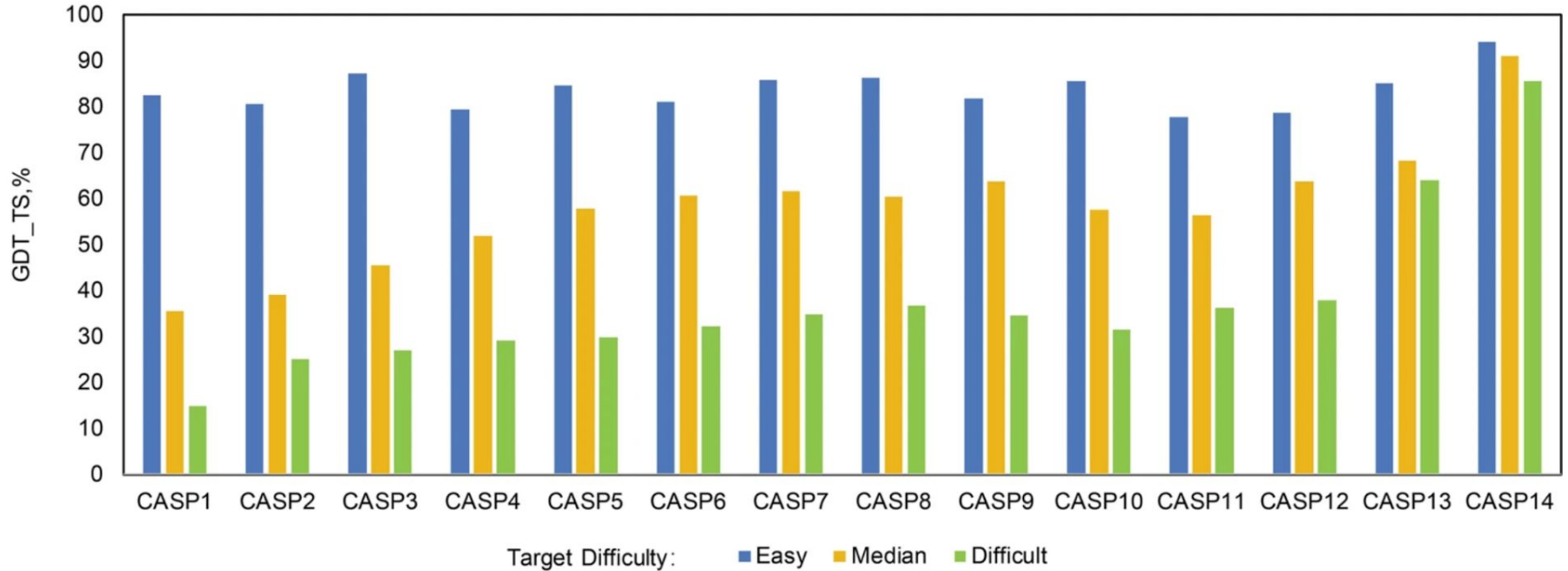*Nature* **420**, 218–223 (2002) | Cite this article

# Computational Methods

# CASP Competition drives protein structure prediction

- Critical Assessment of Structure Prediction (CASP)
- Founded in 1994
- 'Olympics' of protein structure prediction
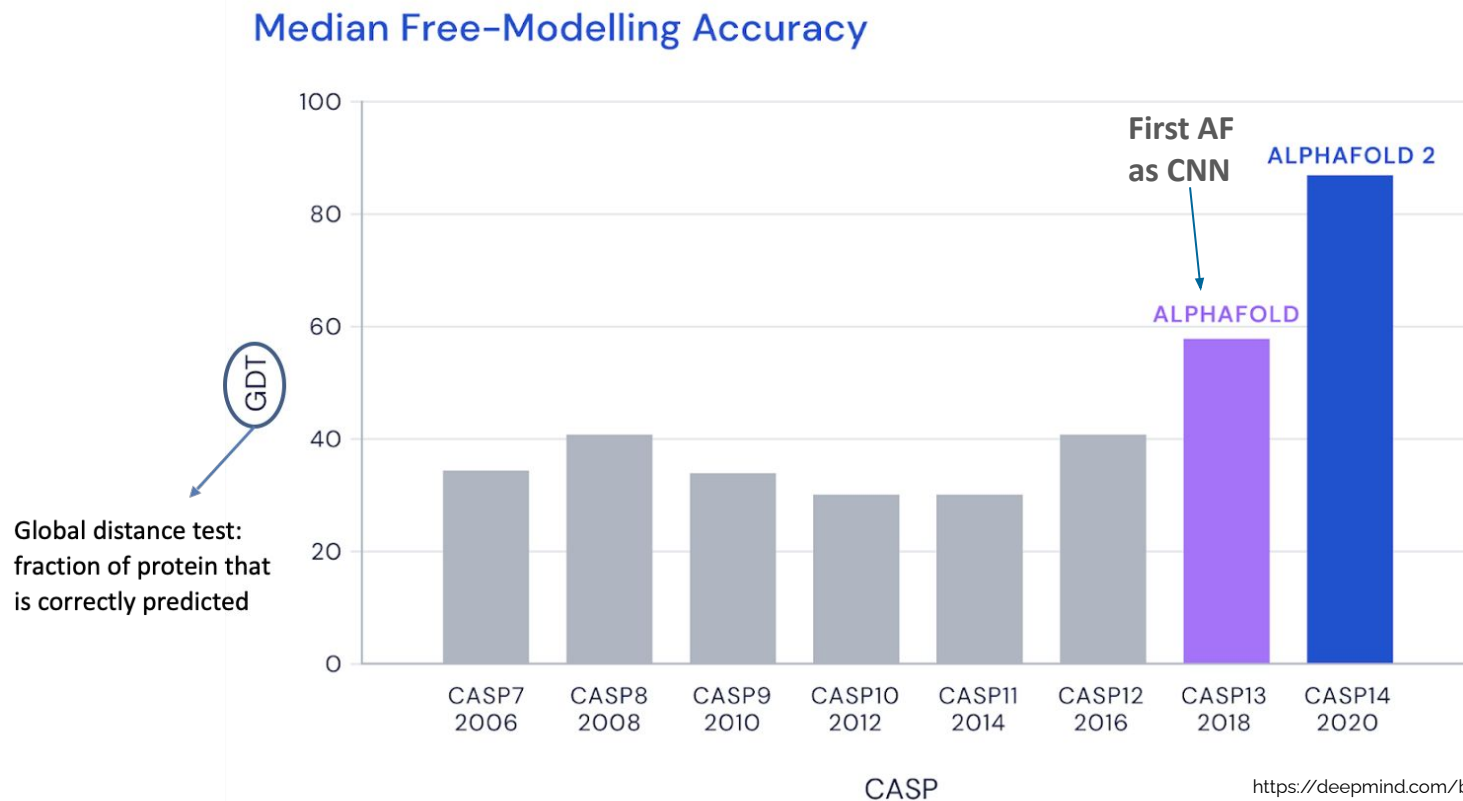- Supported by Lawrence Livermore National Labs, DOE, etc.
- Check it out: https://www.predictioncenter.org

# CASP Competition drives protein structure prediction

# AF2 compared to CASP winners over the years

**Median Free–Modelling Accuracy**

First AF as CNN

ALPHAFOLD 2

ALPHAFOLD

GDT

Global distance test: fraction of protein that is correctly predicted

100
80
60
40
20
0

CASP7 2006
CASP8 2008
CASP9 2010
CASP10 2012
CASP11 2014
CASP12 2016
CASP13 2018
CASP14 2020

CASP

# CASP14 (2020) Results: Entrance of AlphaFold2



| # | GR code | GR name | Domains Count | SUM Zscore (>-2.0) |
|---|---------|---------|---------------|-------------------|
| 1 | 427 | AlphaFold2 | 92 | 244.0217 |
| 2 | 473 | BAKER | 92 | 90.8241 |
| 3 | 403 | BAKER-experimental | 92 | 88.9672 |
| 4 | 480 | FEIG-R2 | 92 | 72.5351 |

# AlphaFold2: the dawn of a new age

# Start Demo!

- [ColabFold Server](ColabFold Server)
- \>2E8I_1|Chain A|6-aminohexanoate-dimer hydrolase|Flavobacterium sp. (37931)
- MNARSTGQHPARYPGAAAGEPTLDSWQEPPHNRWAFAHLGEMVPSAAVSRRPVNAPGHA LARLGAIAAQLPDLEQRLEQTYTDAFLVLRGTEVVAEYYRAGFAPDDRHLLMSVSKSLCGTVV GALVDEGRIDPAQPVTEYVPELAGSVYDGPSVLQVLDMQISIDYNEDYVDPASEVQTHDRSA GWRTRRHGDPADTYEFLTTLRGDGSTGEFQYCSANTDVLAWIVERVTGLRYVEALSTYLWA KLDADRDATITVDTTGFGFAHGGVSCTARDLARVGRMMLDGGVAPGGRVVSEDWVRRVLA GGSHEAMTDKGFTNTFPDGSYTRQWWCTGNERGNVSGIGIHGQNLWLDPLTDSVIVKLSS WPDPDTEHWHRLQNGILLDVSRALDAV
- 392 Amino Acids, Nylon Hydrolase

# What is an MSA?

# MSA: Multiple Sequence Alignment

# Direct-coupling analysis of residue coevolution captures native contacts across many protein families
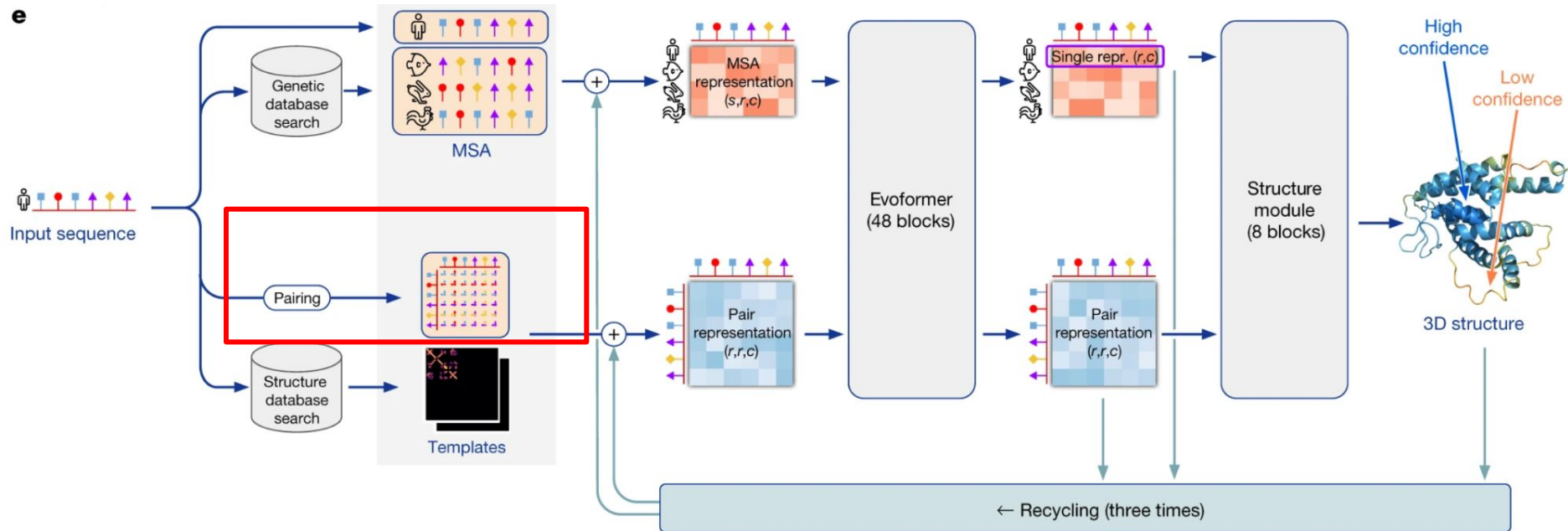


correlated

constraint

inference

contact in 3D

# Initializing the Pair Representation

# What is a pair representation?
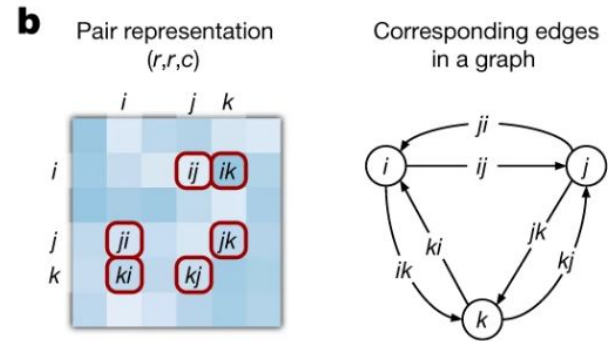


Distogram is used to map 2D pairwise distances
Distogram are independent of translations and rotations, so no need to align structures (much faster)

2D presentation of 3D structure in Distance Map

**b** Pair representation (r,r,c)    Corresponding edges in a graph

Senior, A. W., Evans, R., Jumper, J., Kirkpatrick, J., Sifre, L., Green, T., ... & Hassabis, D. (2020). Nature.

# What is a pair representation?



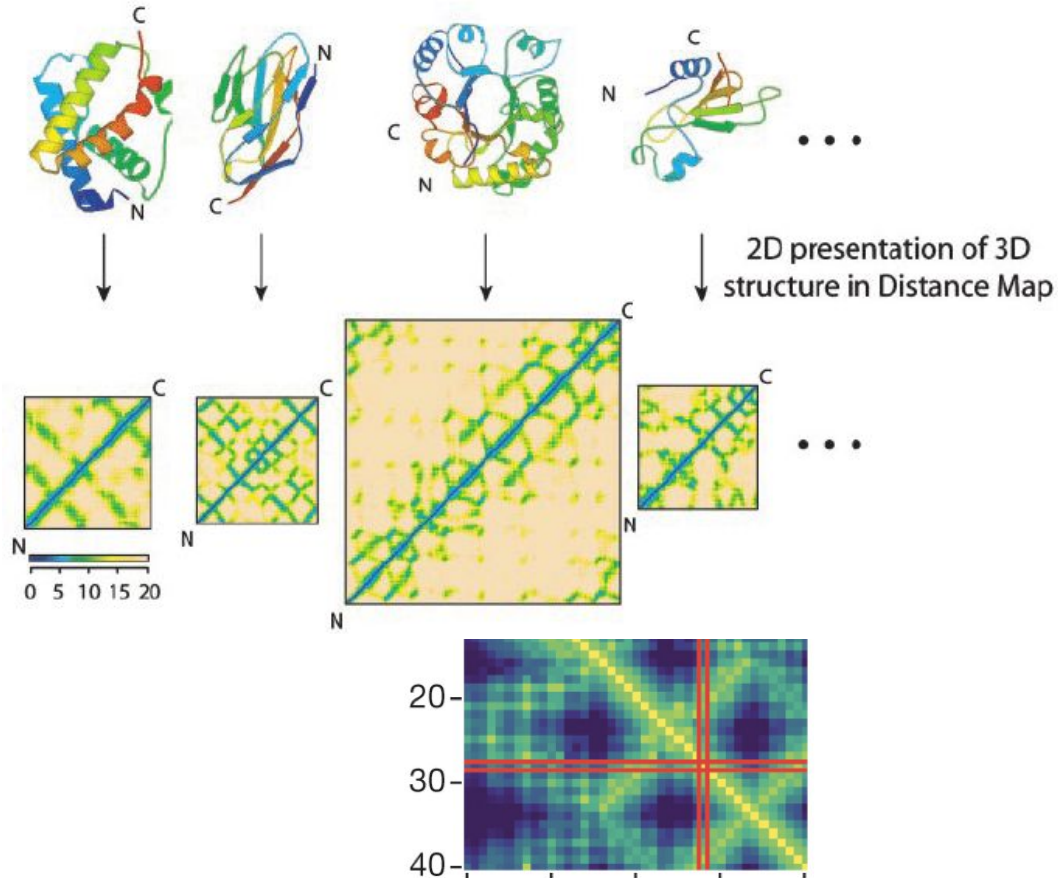2D presentation of 3D structure in Distance Map

0 5 10 15 20

Distogram is used to map 2D pairwise distances
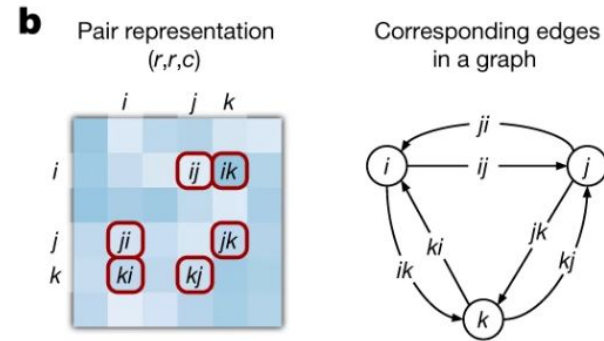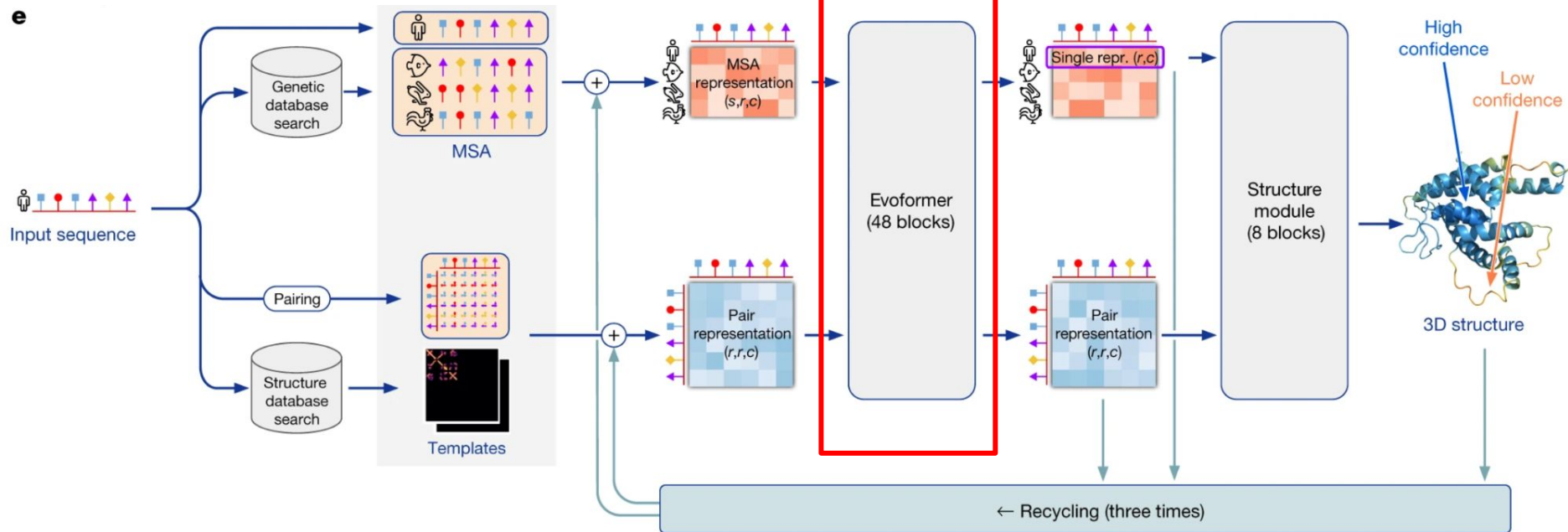Distogram are independent of translations and rotations, so no need to align structures (much faster)

**b** Pair representation ($r,r,c$)

Corresponding edges in a graph

Senior, A. W., Evans, R., Jumper, J., Kirkpatrick, J., Sifre, L., Green, T., ... & Hassabis, D. (2020). Nature.

# Transformer-like module updates representations

# Attention Is All You Need

**Ashish Vaswani**[*]
Google Brain
avaswani@google.com

**Noam Shazeer**[*]
Google Brain
noam@google.com

**Niki Parmar**[*]
Google Research
nikip@google.com

**Jakob Uszkoreit**[*]
Google Research
usz@google.com

**Llion Jones**[*]
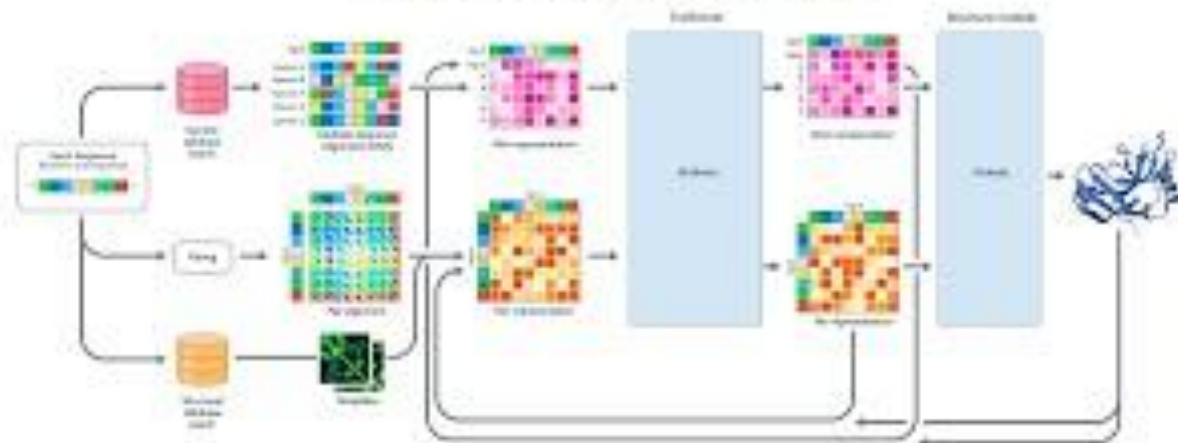Google Research
llion@google.com

**Aidan N. Gomez**[* †]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser**[*]
Google Brain
lukaszkaiser@google.com
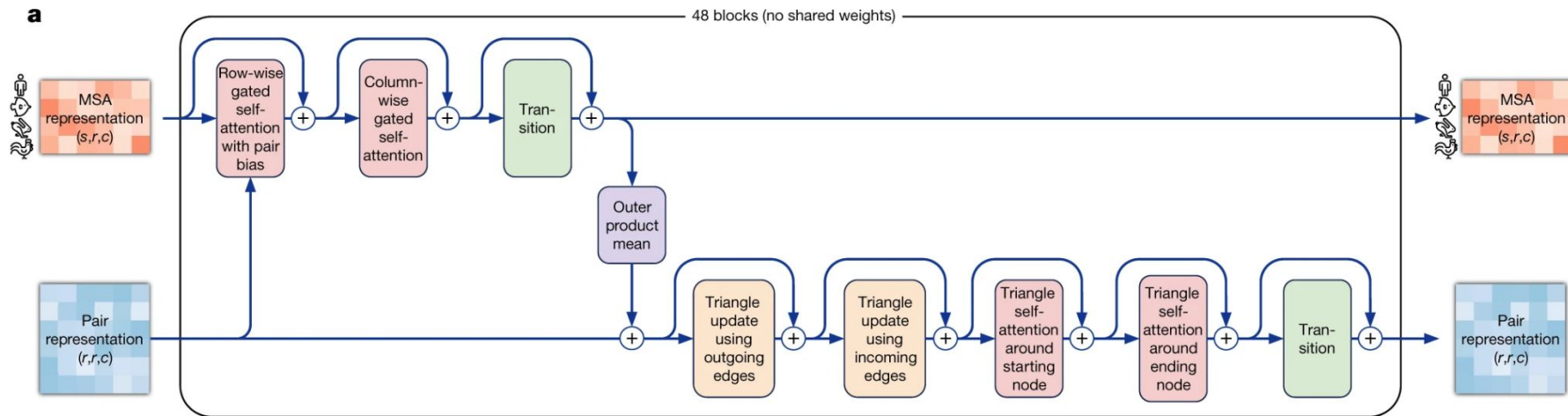
**Illia Polosukhin**[* ‡]
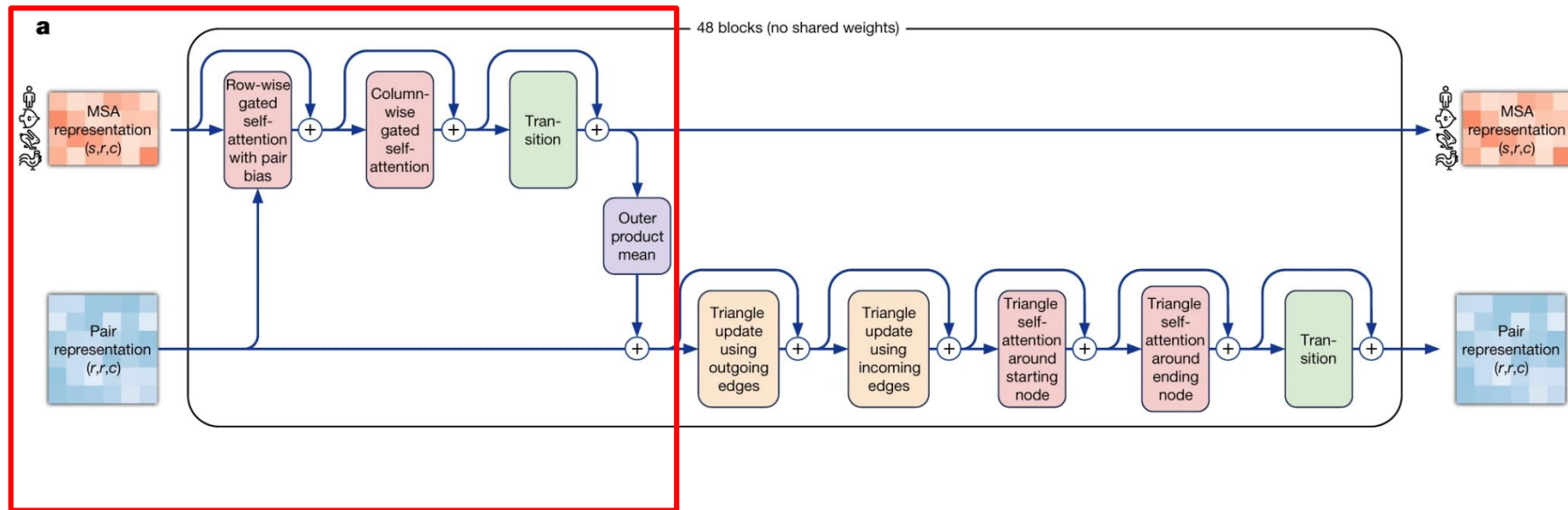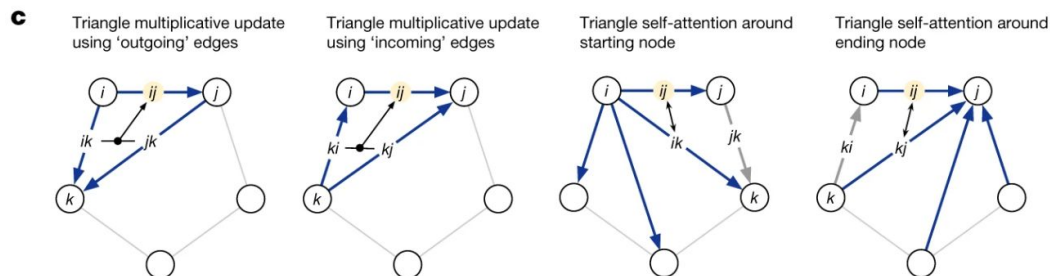illia.polosukhin@gmail.com

What Is AlphaFold?

# The Evoformer: 48 blocks of back and forth between evolutionary and spatial reasoning.
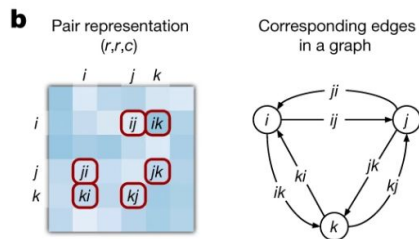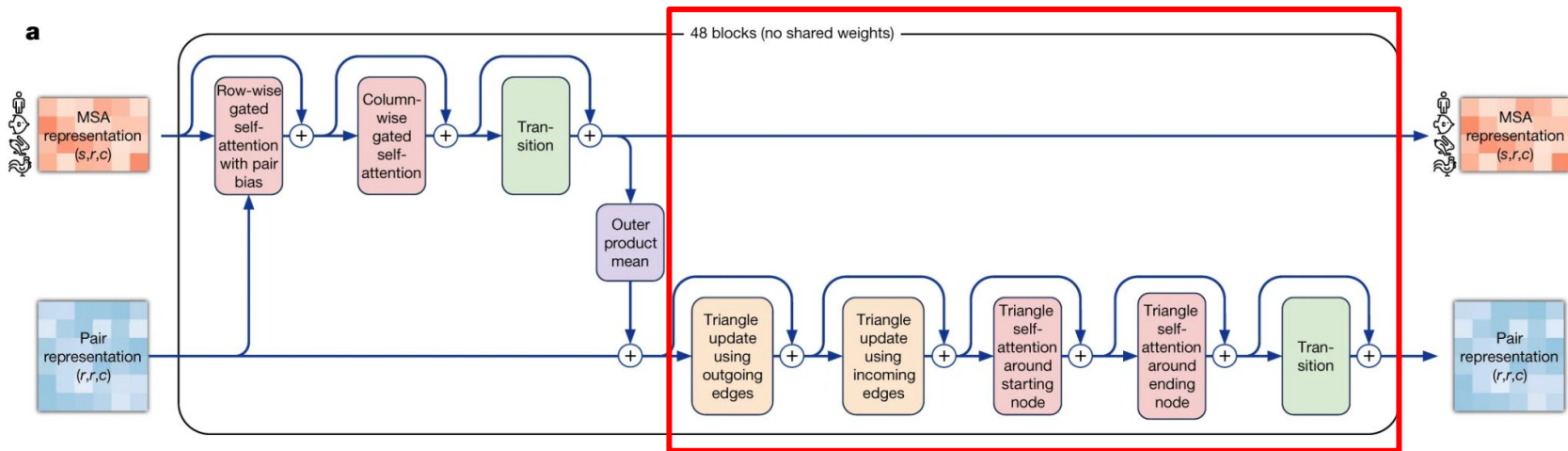
# The Evoformer: 48 blocks of back and forth between evolutionary and spatial reasoning.

# The Evoformer: 48 blocks of back and forth between evolutionary and spatial reasoning.

# Structure module converts representation to structure

# Structure module converts representation to structure

# The Structure Module: iteratively updating a residue gas

# Recycling iteratively refines structure

# How do we know it works? – ask the model!



**T1037 / 6vr4**
90.7 GDT
(RNA polymerase domain)

**T1049 / 6y4f**
93.3 GDT
(adhesin tip)

● Experimental result
● Computational prediction

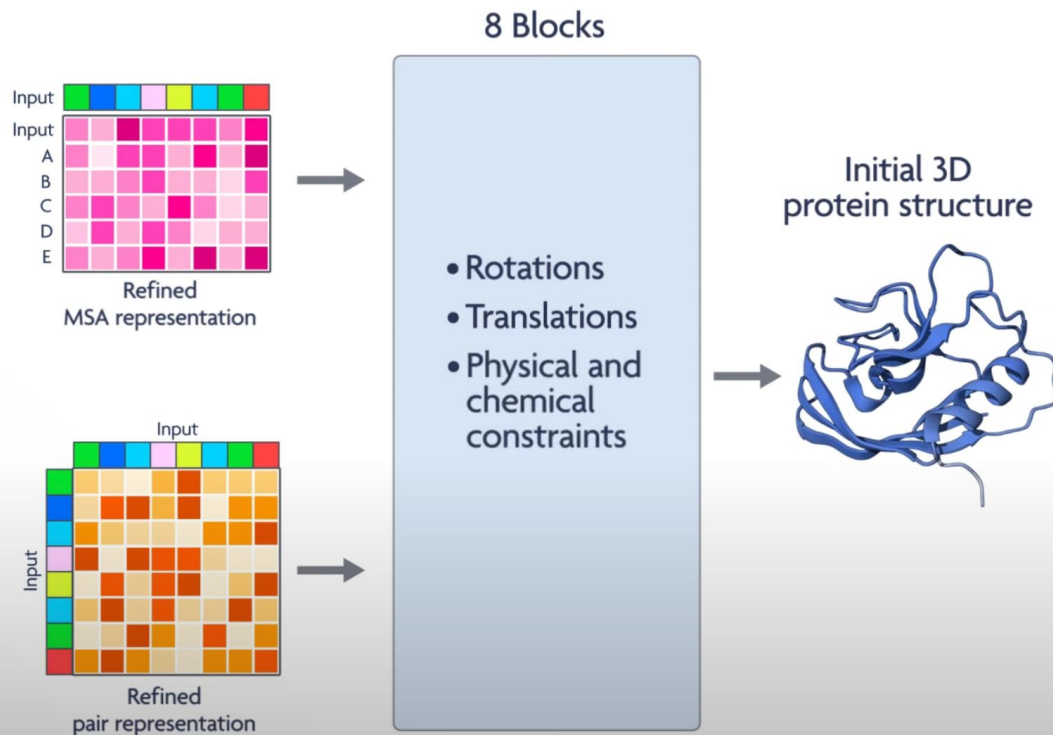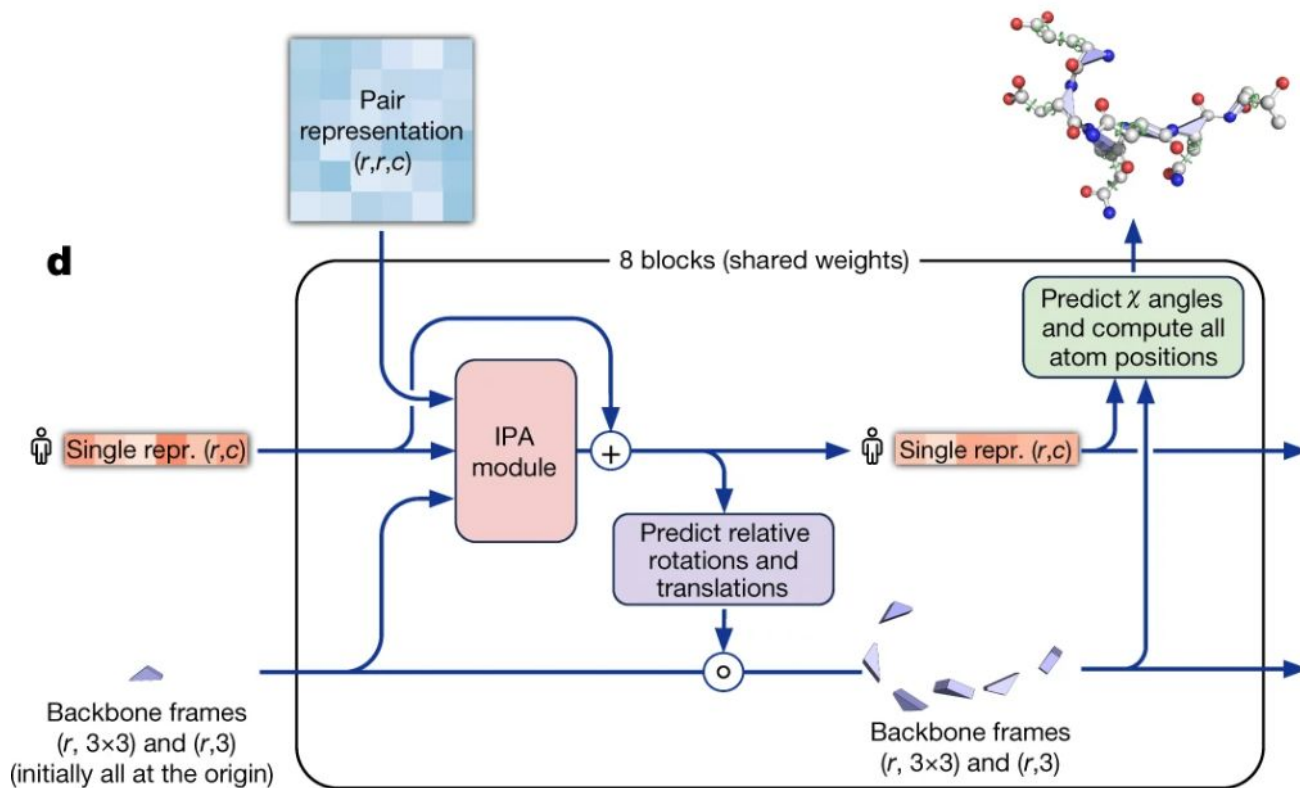# Check in on Demo

# The 3Ps: pLDDT, PAE, and pTM

- LDDT, AE, TM require ground truth. What can we do?

- pLDDT: predicted **L**ocal **D**istance **D**ifference **T**est

- PAE: **P**redicted **A**ligned **E**rror

- pTM: predicted **T**emplate **M**odeling score
  - Global comparison of similarity between two structures
  - Measure of 0 to 1

# The 3Ps: pLDDT, PAE, and pTM

- LDDT, AE, TM require ground truth. What can we do?

- pLDDT: predicted **L**ocal **D**istance **D**ifference **T**est

- PAE: **P**redicted **A**ligned **E**rror

- pTM: predicted **T**emplate **M**odeling score
  - Global comparison of similarity between two structures
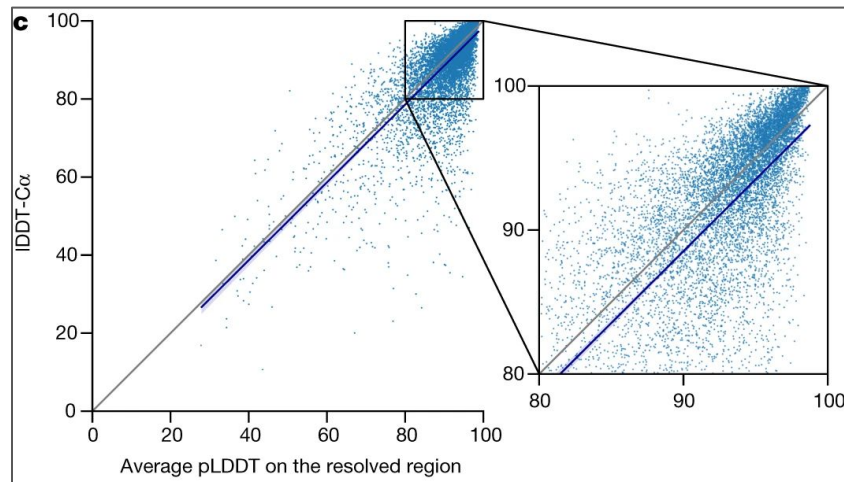  - Measure of 0 to 1
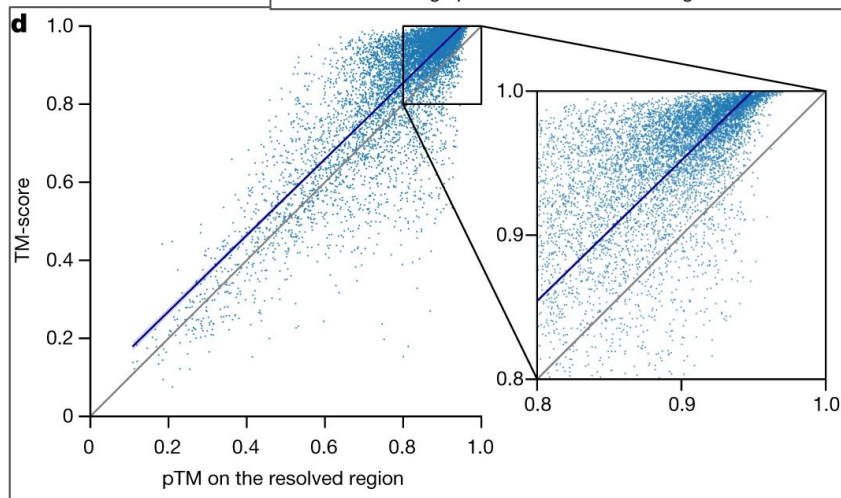
# The 3Ps: pLDDT, PAE, and pTM

- LDDT, AE, TM require ground truth. What can we do?

- pLDDT: predicted **L**ocal **D**istance **D**ifference **T**est

- PAE: **P**redicted **A**ligned **E**rror

- pTM: predicted **T**emplate **M**odeling score
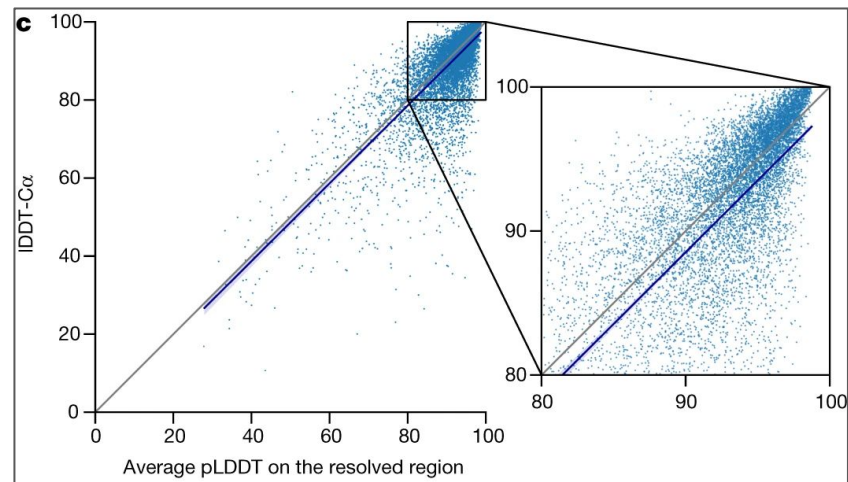  - Global comparison of similarity between two structures
  - Measure of 0 to 1



Highly accurate protein structure prediction with AlphaFold. *Nature* 596, 583–589 (2021). https://doi.org/10.1038/s41586-021-03819-2

# Predicted Local Distance Difference Test

- AlphaFold's per-residue prediction of its lDDT-Ca score
- Low lDDT commonly associated with disorder
- High pLDDT on each domain doesn't imply confidence in relative positions!



Per residue confidence

High (> 90)   Low (50 – 70)
Medium (70 – 90)   Very low (< 50)

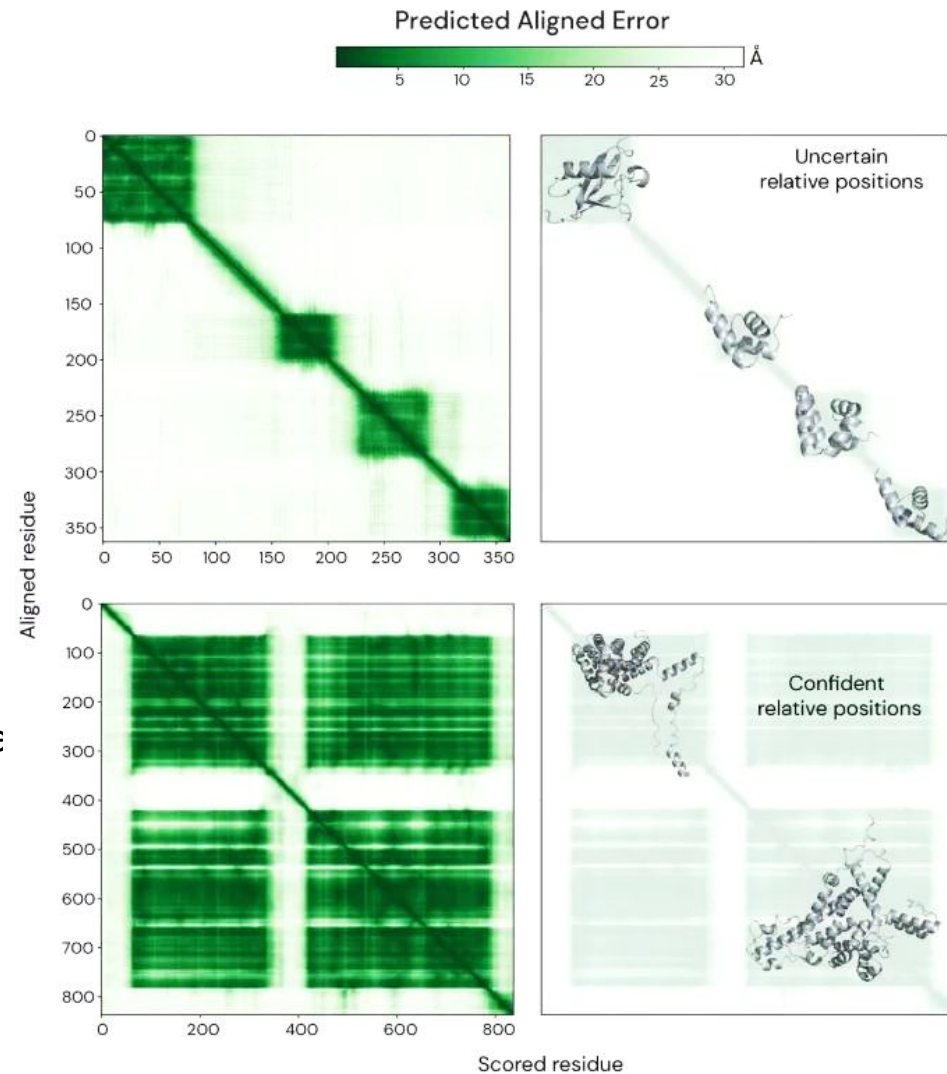# Predicted Aligned Error (PAE)



- Prediction of position error at residue x if the predicted and the true structures were aligned on y
- PAE aims to measure confidence in the relative positions of pairs of residues
- Use where pairwise confidence is relevant – interpreting domain distances in a multi domain protein
- Suppose residue y were aligned to the true structure and we measured the position error at residue x. The color at (x,y) is AF's prediction of that error
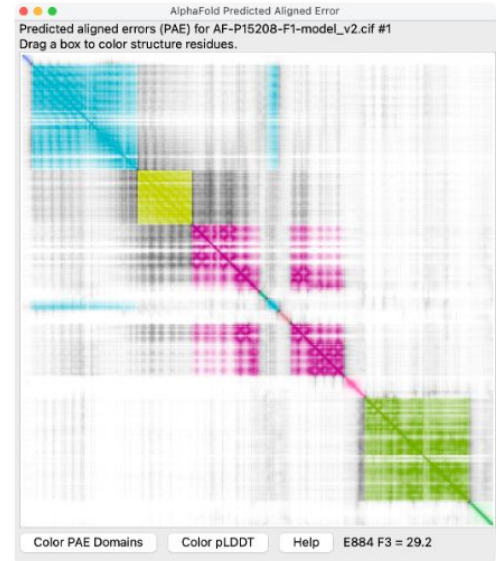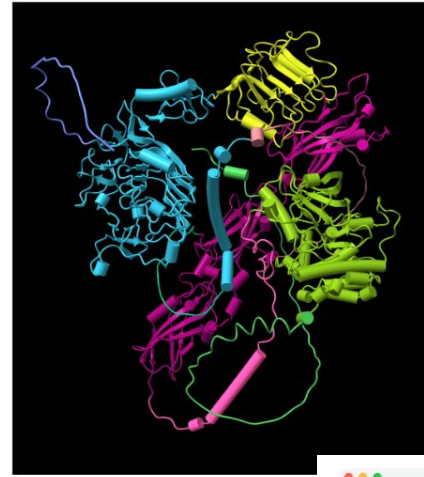
# Predicted Aligned Error (PAE)



- Prediction of position error at residue x if the predicted and the true structures were aligned on y
- PAE aims to measure confidence in the relative positions of pairs of residues
- Use where pairwise confidence is relevant – interpreting domain distances in a multi domain protein
- Suppose residue y were aligned to the true structure and we measured the position error at residue x. The color at (x,y) is AF's prediction of that error

# 214,683,839 Predicted Protein Structures on AFDB!

Proteo...



N

pLDDT ∈ [90–100]
pLDDT ∈ [70–90)
pLDDT ∈ [50–70)
pLDDT ∈ [0–50)
PDB 6YJ1

**ChimeraX**
@UCSFChimeraX

The human muscle protein titin predicted by AlphaFold, 34350 residues.

7:23 PM · Sep 8, 2021 · Twitter Web App

...sically ...ered

... 20% ...y

... 95% ...y

Dark proteome
Intrinsically disordered
PFAM – No PDB / AF
Gained AF – 70 > pLDDT ≥ 50
Gained AF – 90 > pLDDT ≥ 70
Gained AF – pLDDT ≥ 90
PDB 20% to 50% - pLDDT < 90
PDB 50% to 95%
PDB ≥ 95%

Tunyasuvunakool, K., Adler, J., Wu, ...

..., Valentini, S., & Valencia, A.
...iology.

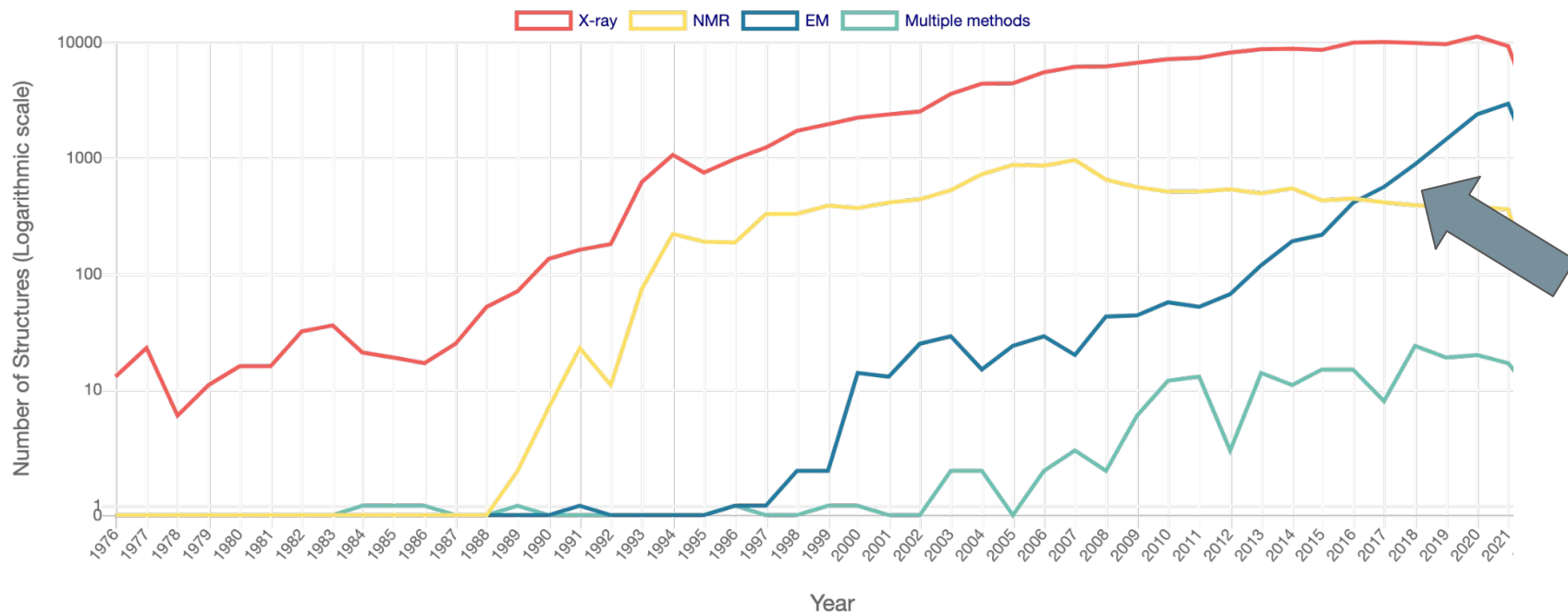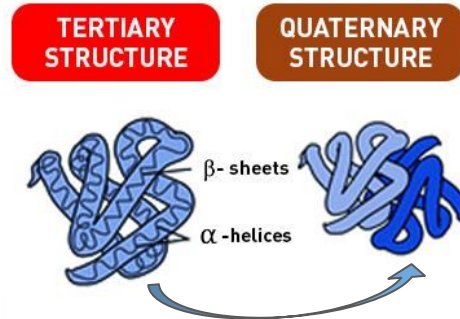# 2024 Nobel Prize in Chemistry



David Baker, Demis Hassabis and John Jumper (left to right) won the chemistry Nobel for developing computational tools that can predict and design protein structures.   Credit: BBVA Foundation
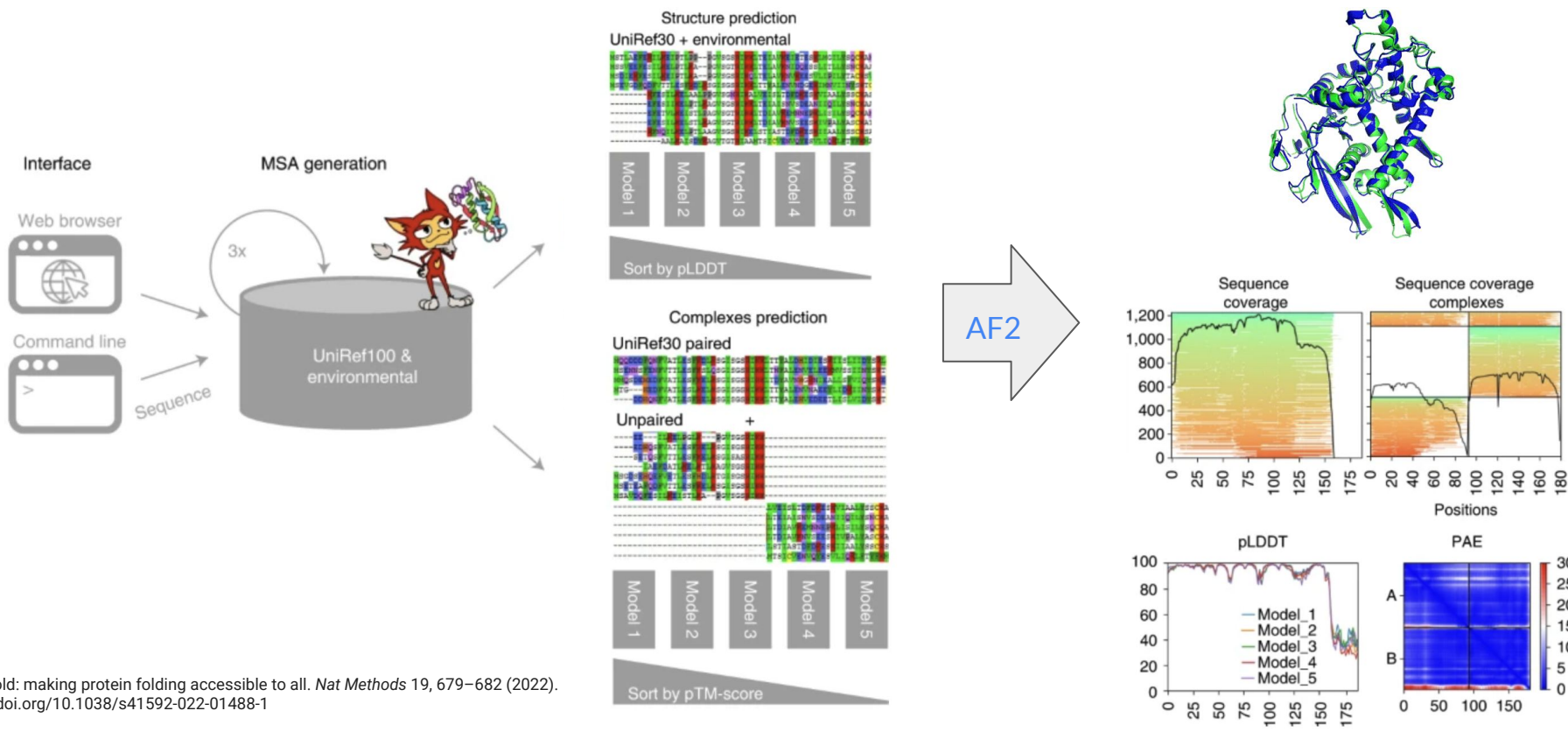
# AlphaFold Multimer!

# ColabFold speeds up AlphaFold by 40-60 fold!

ColabFold: making protein folding accessible to all. *Nat Methods* 19, 679–682 (2022).
https://doi.org/10.1038/s41592-022-01488-1
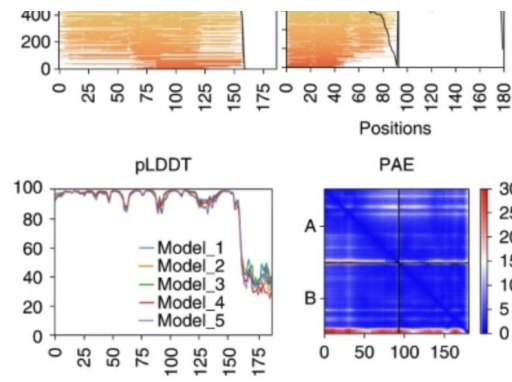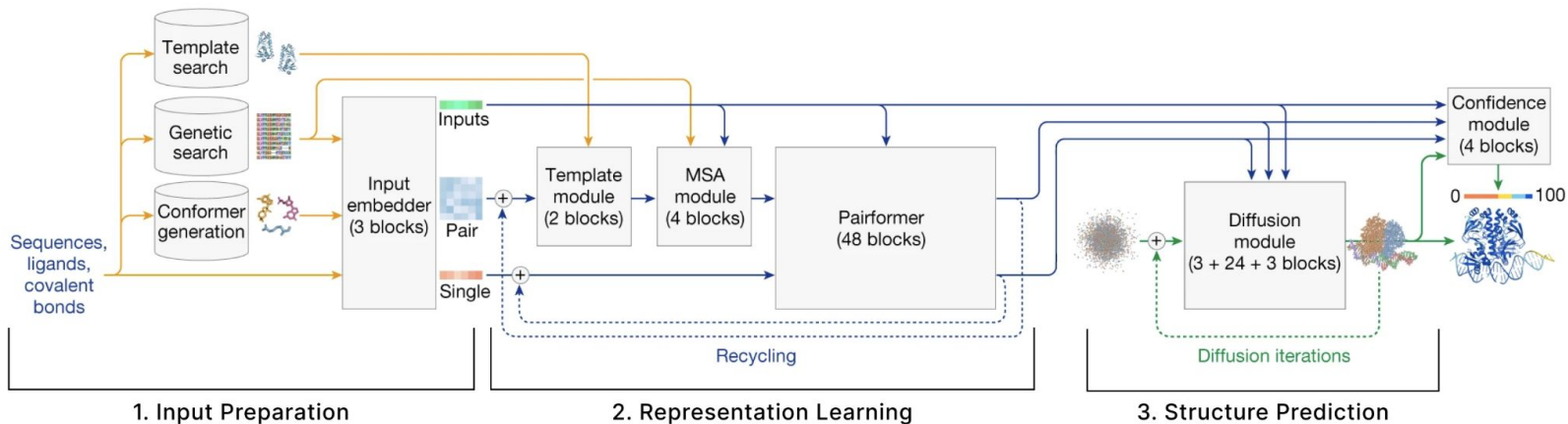
# ColabFold speeds up AlphaFold by 40-60 fold!



"ColabFold's 40−60-fold faster search and optimized model utilization enables prediction of close to 1,000 structures per day on a server with one graphics processing unit."

# A new era: All Atom Models (AF3)

**Accurate structure prediction of biomolecular interactions with AlphaFold 3**

Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J. Ballard, Joshua Bambrick, Sebastian W. Bodenstein, David A. Evans, Chia-Chun Hung, Michael O'Neill, David Reiman, Kathryn Tunyasuvunakool, Zachary Wu, Akvilė Žemgulytė, Eirini Arvaniti, Charles Beattie, Ottavia Bertolli, Alex Bridgland, Alexey Cherepanov, Miles Congreve, … John M. Jumper ✉    + Show authors

# The Illustrated AlphaFold

A visual walkthrough of the AlphaFold3 architecture, with more details and diagrams than you were probably looking for.

AUTHORS
Elana Simon
Jake Silberg

AFFILIATIONS
Stanford University
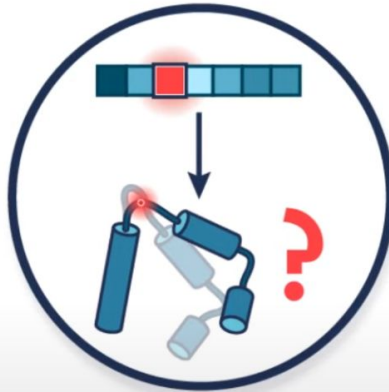Stanford University

PUBLISHED
July 10, 2024

# Applications and Frontiers!



Discovering drugs

Effect of genetic variants

Modeling protein–protein interactions

Engineering artificial proteins

# Has protein structure prediction been solved?

- **Sort of.**
- **The Protein Folding Problem (de novo)**
  - What is the folding code?
  - **What is the folding mechanism?**
  - Can we predict a native protein structure from its primary, amino acid sequence?
    - **No for a sequence in isolation…**
    - **Yes when informed by like sequences and their structures**

# Resources & Useful Links

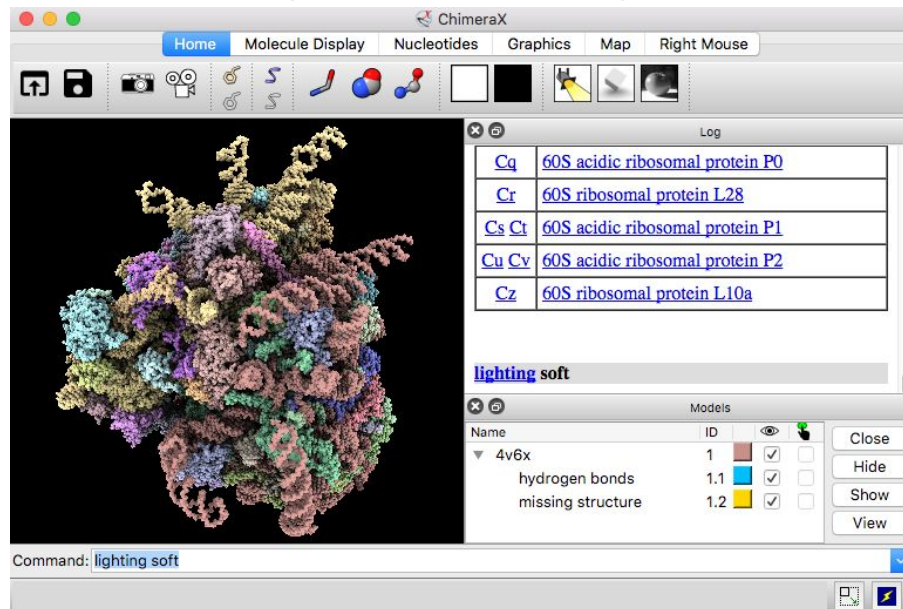- Prediction Servers/Colabs:
  - AF3 Server
  - ColabFold (AF2 w/MMSeqs2)
  - Maintained list of Google Colabs: https://github.com/sokrypton/ColabFold
- AlphaFold Resources:
  - Github page: https://github.com/google-deepmind/alphafold
  - AFDB Protein Structure Database and links: https://alphafold.com
  - The Illustrated AlphaFold
  - Lovely & more in depth lecture from John Jumper at Vanderbilt University
  - Running AlphaFold on the BRCF Servers:
    - AMD GPUs
    - NVIDIA GPUs
  - Running AlphaFold on TACC
  - AlphaFold2 & Equivariance

# ChimeraX

- Download ChimeraX: https://www.cgl.ucsf.edu/chimerax/download.html
- Quick Start: https://www.cgl.ucsf.edu/chimerax/docs/quickstart/index.html
- Very comprehensive user guide: https://www.cgl.ucsf.edu/chimerax/docs/user/index.html

# Thank you! Questions?

Looking for people to work with/learn more about Machine Learning applied to Biology?

https://www.biomlsociety.org

We meet every other Thursday from 11am to noon, in MBB 3.204

**TACOS** and **COFFEE** provided!

Watch the Commander Complex assemble!