DeepMind

# Highly accurate protein structure prediction with AlphaFold

**Tim Green**

*Learning to Discover, Paris*

John Jumper[1]*†, Richard Evans[1]*, Alexander Pritzel[1]*, Tim Green[1]*, Michael Figurnov[1]*, Kathryn Tunyasuvunakool[1]*, Olaf Ronneberger[1]*, Russ Bates[1]*, Augustin Žídek[1]*, Alex Bridgland[1]*, Clemens Meyer[1]*, Simon A A Kohl[1]*, Anna Potapenko[1]*, Andrew J Ballard[1]*, Andrew Cowie[1]*, Bernardino Romera-Paredes[1]*, Stanislav Nikolov[1]*, Rishub Jain[1]*, Jonas Adler[1], Trevor Back[1], Stig Petersen[1], David Reiman[1], Martin Steinegger[2], Michalina Pacholska[1], David Silver[1], Oriol Vinyals[1], Andrew W Senior[1], Koray Kavukcuoglu[1], Pushmeet Kohli[1], Demis Hassabis[1]*†

[1]DeepMind, London, UK, [2]Seoul National University, South Korea
* Equal contribution † Corresponding authors: John Jumper (jumper@deepmind.com), Demis Hassabis (dhcontact@deepmind.com)
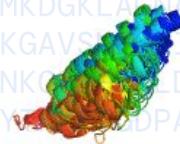
DeepMind

# Introduction

# About me

Tim Green – tfgg@deepmind.com

- 2006–2010 MSci at Cambridge Physics

- 2010–2014 DPhil at Oxford Materials (MML)

- 2014–2015 Postdoc at same

- 2016–present DeepMind

Thesis "Prediction of NMR J-coupling in condensed matter" developed DFT predictions with relativistic effects, disorder and temperature.

Work at DeepMind: *Computer Vision, Population Based Training, Deep RL, ML infrastructure* and *protein structure prediction*

# DeepMind and protein folding

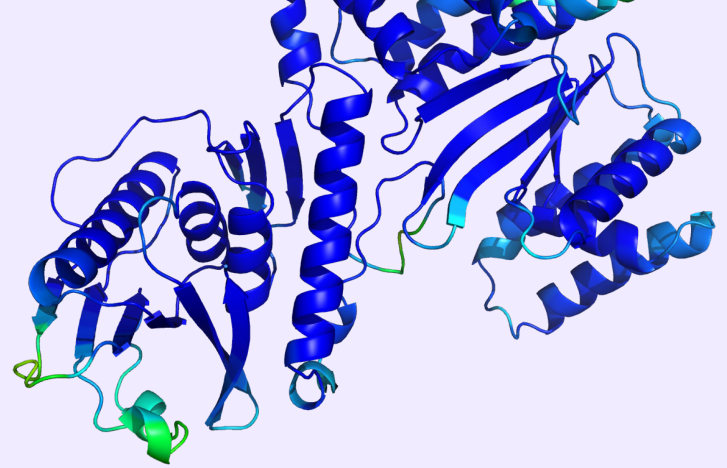**A central part of DeepMind's mission is to solve fundamental scientific problems with AI**

Predicting the 3D structure of a protein from its amino acid sequence is one such challenge

AlphaFold is our model that aims to solve this problem

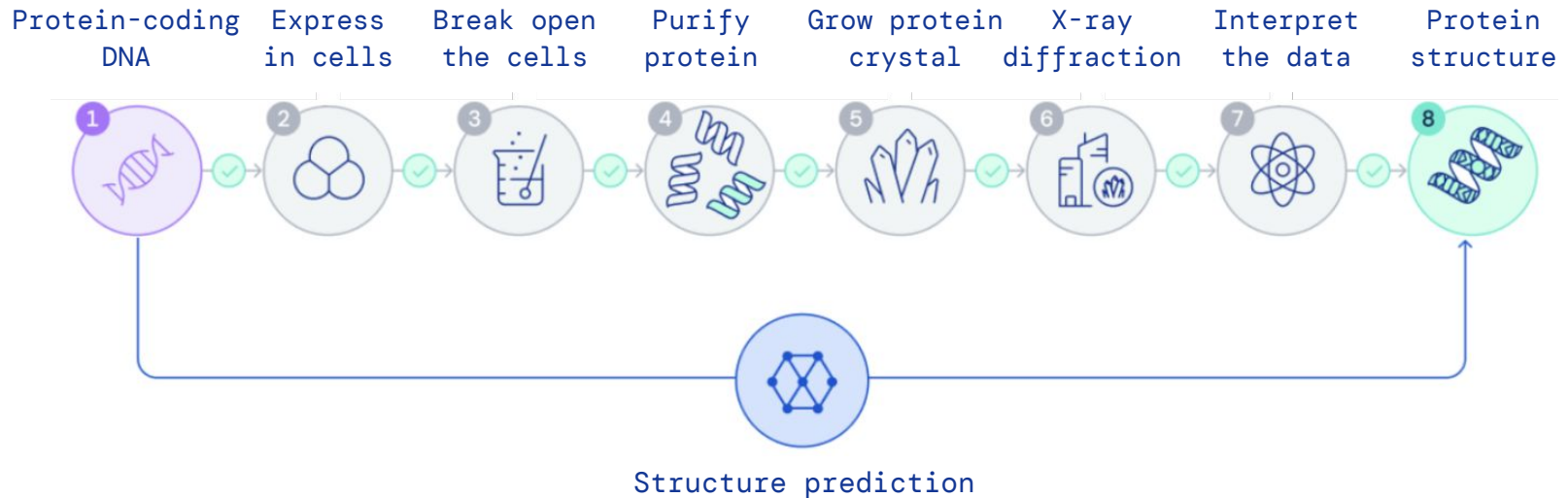IAKFGDWINDEVERNVNEDGEPLLIQDVRQDSSKHYFFILKNGER
FDLLTREFDSFTSPDLTNEIKEITDQLSYYIYNKHFSSDFEQVEG
AKLNIQNEISQFVKEGKAPVQAAYNKLQDPDIKDLLDYYDNIEKH
SDEFESEIVKFFSEKKLIIKDAELEDVTQEGLNEGLQGGDLVQAF
EKNSKDNATANVKLMLSFLPKIDNLTGEPALGDYLNKPVFRSFDS
IHSELLEVLSDITTLHVQGEVLDVFSSMYNKIKELADFKKSFKPL
LEILDTIDEQKKTEFVQAFYLSKINFYTTTIETLETEDQNNTLTT
FKVQNVSNANNPISSKLTEYYTNFKYKILPGGKLNKGKLKDLQST
VTSLLEKTRKENNPKYKSDSDFYEVFEEGVVELMQVFEDLGVDSI
TFEAMDIFLKQFRFDLPENNAYKIMYQQYQGKLTNLNNLLKDIQS
NKINPYKINPFKNYSNLIFNSLAEAENYFIENNNESTIFSNGKTY
WNFARPSYISNRINTFKNNPGVLRQLLNTSYGQSSLWAKHLLGEE
KNVTGDFVLAGNARESASENRLKSLELSIFNSLQEKDKGAEGNDN
GSISIVDQLADKLNKVLRGGTKNGTSIYSTVTPGDKSTLHEIKID
HFIPETISSFSNGTMIFNDKIVNAFTDHFVSEVNRMKEAYQELET
LPESKRVVHYHTDARGNVMKDGKLAGNAFKSGHILSELSFDQITQ
DDNEMLKLYNEDGSPINPKGAVS...ILIKQTINKVLNQRIKEN
IRYFKDQGLVIDTVNKDGNK...DKSIMSEYTDDIQLTEFD
ISHVVSDFTLNSILASIEY...DPANYKNMVDFFKRVPATYT
NGTNLRLGLEANDHLFDVAVLENIVKPSAYLKEIGESLKLSDLSE
AEKKYILEAYEDVNQTDAQAWITPKRWAFLISRTGKWNSKYQSVY
NKILKSESLDASEMKLAAQPLKGVYFGLVNNTPTYLKYSQAVLLP
QLVAGTQLQSLADAMNKQDIGESIVLDGVKVGATTPNIVTDENGD
ILKSISLNPLTLSNADWKLQQDLPVKTIKPTLLGSQIQKNIYSSL
TDEATYTIENEAFNGSGMFQAINDTVSAMSNLSIAGLSSE...DS
EGKIDKRKLYDMLEREMLDKGSAINLLKSIQKNLPIEAMP...R
LYNIVFSKINSAAVKLKTNGGSFIQLSNFGLDKQTADAKGITWLV

# What are proteins?

○ Proteins are molecular machines that are **essential to life**

○ They have **many functions**: from our hair to our immune system

○ Consist of **chains of amino acids** that fold into a 3D structure

○ The exact **3D shape** is important for a protein's function

○ Understanding protein structures is a **fundamental problem in biology**
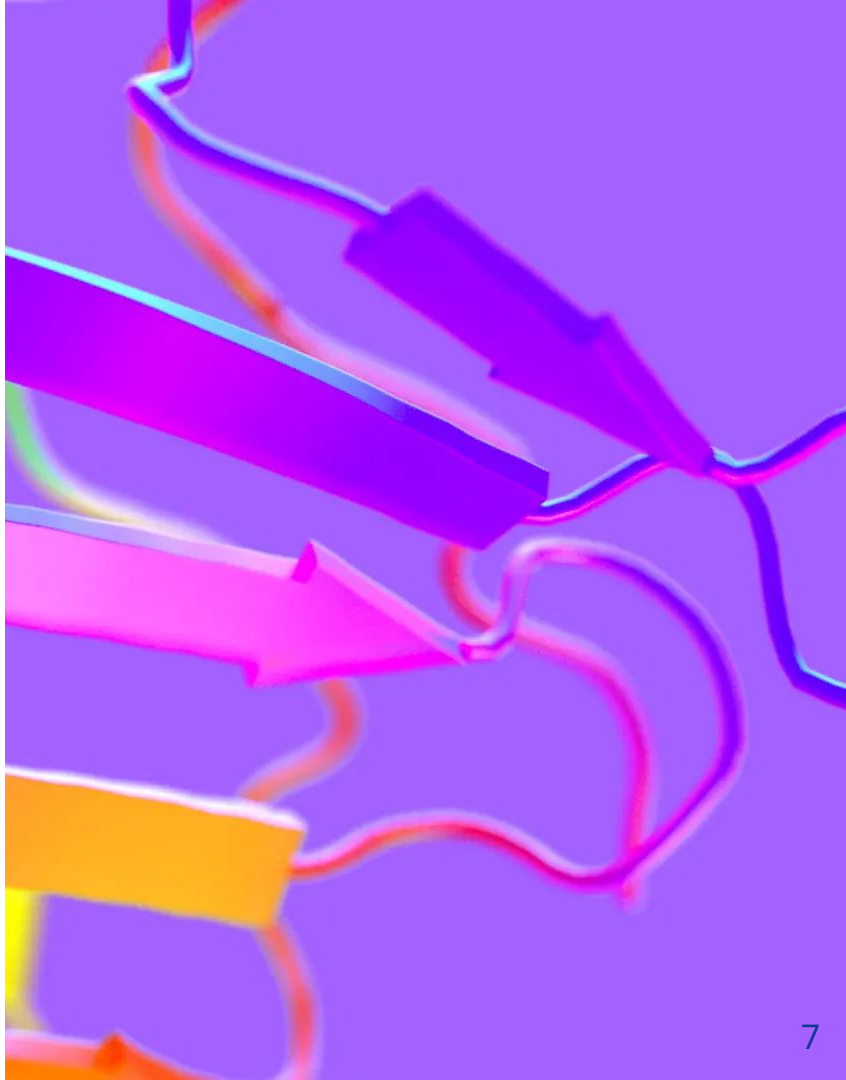
# Why predict protein structures?

Experimental structure determination takes months to years.

Structure prediction can provide actionable information faster.



Protein-coding DNA → Express in cells → Break open the cells → Purify protein → Grow protein crystal → X-ray diffraction → Interpret the data → Protein structure

Structure prediction

# Agenda

- Introduction

- AlphaFold and CASP

- How AlphaFold works

- How AlphaFold builds protein structures

- AlphaFold impact

- AlphaFold–Multimer

DeepMind

# AlphaFold and CASP

# AlphaFold at CASP

When working on these problems, a clear success metric is crucial

Fortunately, the protein structure prediction community had established CASP

The CASP assessment involves predicting recently solved structures that aren't yet public

At CASP14, AlphaFold was the top ranked method achieving consistently high accuracy



9

# CASP: historical perspective



- CASP has provided a biennial blind assessment of structure prediction methods over the last 25 years

- AlphaFold 2 achieved a **median accuracy of 92.4 GDT** over all targets in CASP14

- In response, **AlphaFold was recognised as a solution to the structure prediction problem** by the CASP organizers

# Protein example: T1064 (ORF8)



**T1064 / 7jtl**
87.0 GDT
(ORF8, SARS-CoV-2)

**Ground truth**
**Prediction**

*7JTL: Flower, T.G., et al. (2020) Structure of SARS-CoV-2 ORF8, a rapidly evolving coronavirus protein implicated in immune evasion. Biorxiv.*

# Protein example: T1044 (RNA Polymerase)



→ Folding as a single long chain

**Individual domains**

T1041          T1042          T1043

*6VR4: Drobysheva, A.V., et al. Structure and function of virion RNA polymerase of a crAss-like phage. Nature (2020). (CASP14 target T1044)*

**Ground truth**
**Prediction**

DeepMind

# How AlphaFold works

# Determining Structure from Evolution - Intuition



**Sequence**

**Physics**

**Structure**

**Evolutionary Structure Prediction**

Evolutionary History

**PPM**

**Function**

**Mutate & Select**

# Model Inputs and Outputs

## Inputs

◆ Amino acid sequence (residues) for the protein

◆ Evolutionary-related sequences (sequences that fold to the same structure, but whose amino acid sequence has diverged due to mutations)

## Training Data

◆ **Labelled Data** – 170k structures, 40k after deduplicating

◆ **Unlabelled Data** – 350k deduplicated sequences

## Output

◆ 3D position of every atom in the protein (300 – 50,000 atoms)

SQETRKKCTEMKKKFKNCEVRCDESNHCVEVRCSDTKYTLC

Multiple Sequence Alignment (MSA)

# Inductive Bias for Deep Learning Models



**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)

# Putting our protein knowledge into the model

→ Physical and geometric insights are built into the network structure, not just a process around it

→ End-to-end system directly producing a structure instead of inter-residue distances

→ Inductive biases reflect our knowledge of protein physics and geometry
  ○ The positions of residues in the sequence are de-emphasized
  ○ Instead residues that are close in the folded protein need to communicate
  ○ The network iteratively learns a graph of which residues are close, while reasoning over this implicit graph as it is being built

# Network



Feeding certain outputs back through the network again improves performance

# Structure module

→ **End–to–end folding** instead of gradient descent

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

→ **3–D equivariant transformer architecture** updates the rigid bodies / backbone
  ○ Also builds the side chains from torsion angles



*Image: Dcrjsr, vectorised Adam Rędzikowski (CC BY 3.0, Wikipedia)*

Iteration 1

**Target: T1041**

19

# Using unlabelled sequences in training

We know ~200k protein structures (Protein Data Bank) but several billion protein sequences.

We use these data in two ways.



MSA BERT

(train model to predict masked locations in MSA)

Noisy student self-distillation

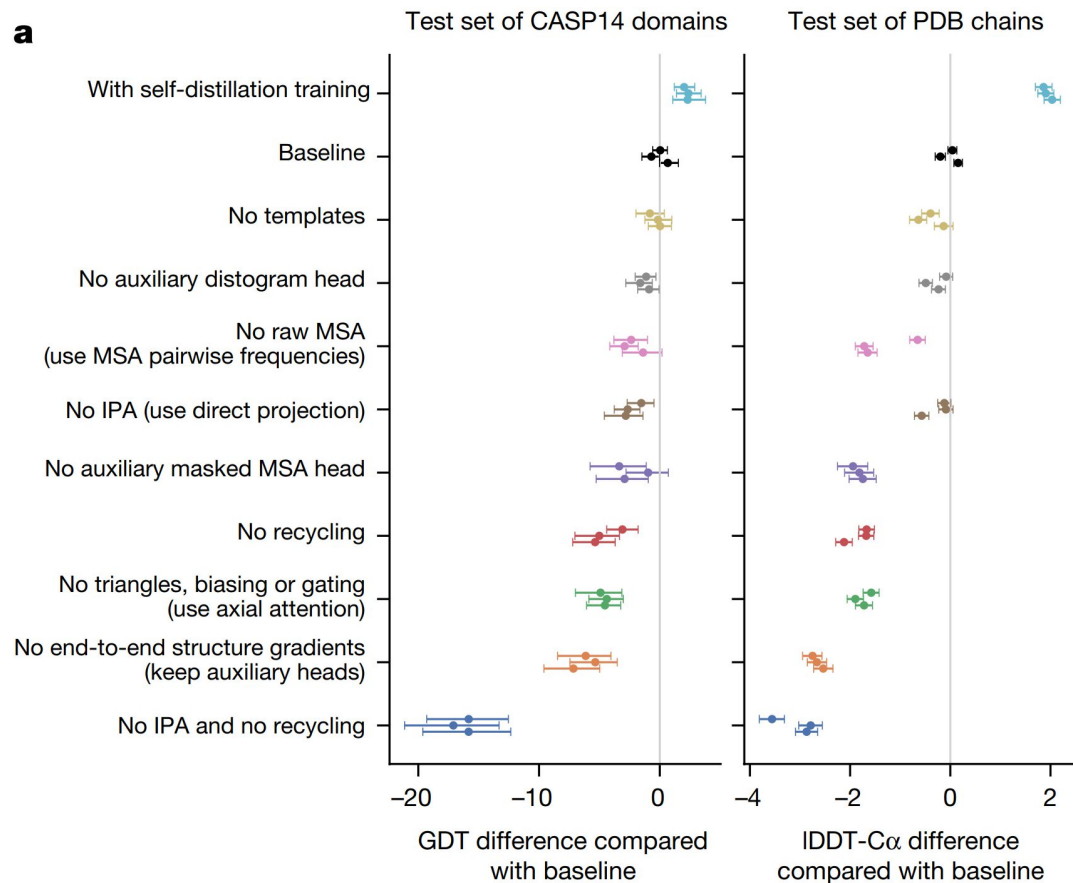(train from predictions of same architecture)

# Which parts mattered? All of it

**No single improvement is dominant**

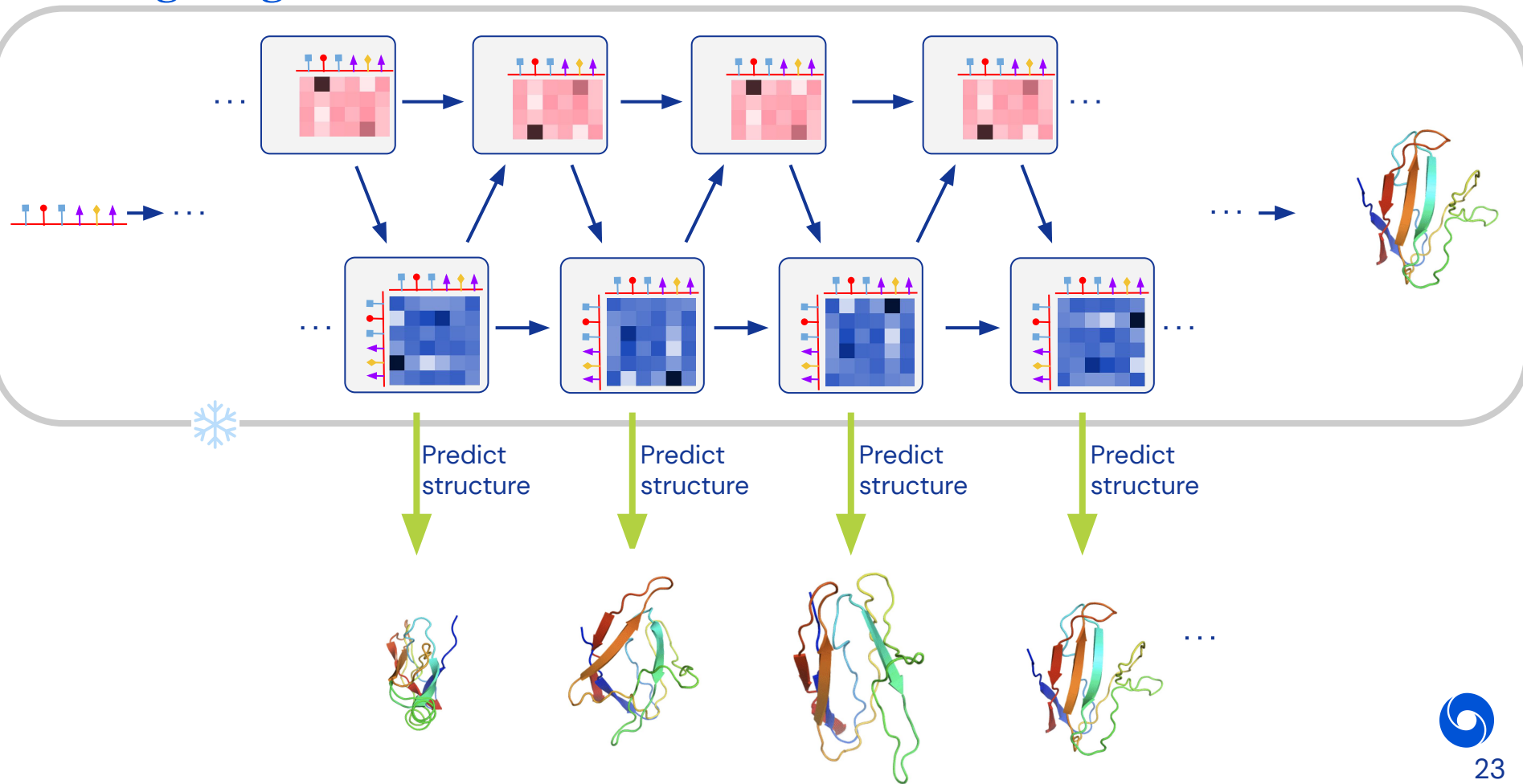More important was the methodology of building the protein intuition into the model

Multiple ablations suggest strong interactions between many of the components



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

lDDT-Cα difference compared with baseline

DeepMind

# How AlphaFold builds protein structures

# Interrogating the Network



Predict structure

Predict structure

Predict structure

Predict structure

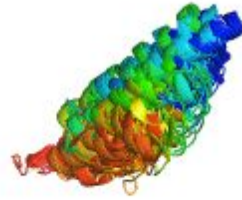# Model interpretability - ORF8 - Sars-Cov2



7JTL: Flower, T.G., et al. (2020) Structure of SARS-CoV-2 ORF8, a rapidly
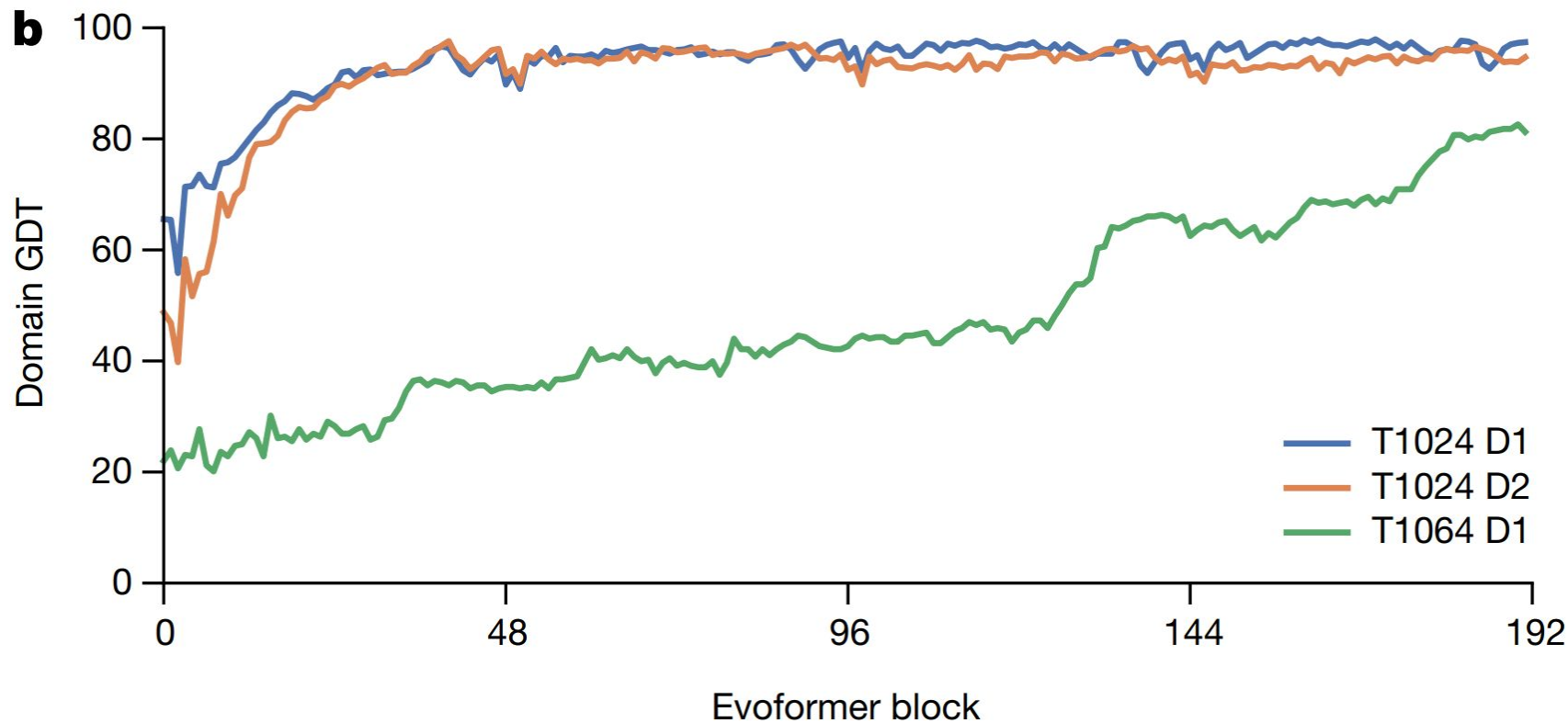evolving coronavirus protein implicated in immune evasion. Biorxiv.

# Model interpretability - T1044

# Model interpretability - Role of depth
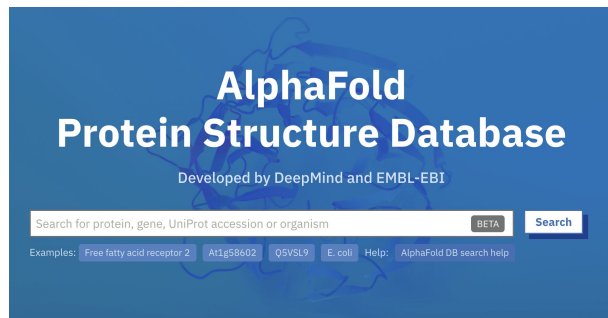
DeepMind

# AlphaFold Impact

# Open Source & AlphaFold Protein Structure Database

We open sourced the code and model weights to run AlphaFold – github.com/deepmind/alphafold (8.3k ⭐)

Also created AlphaFold Protein Structure Database

- Website developed and hosted by EMBL–EBI

- Contains pre-run predictions for **21 model organisms + SwissProt** (>800k structures)

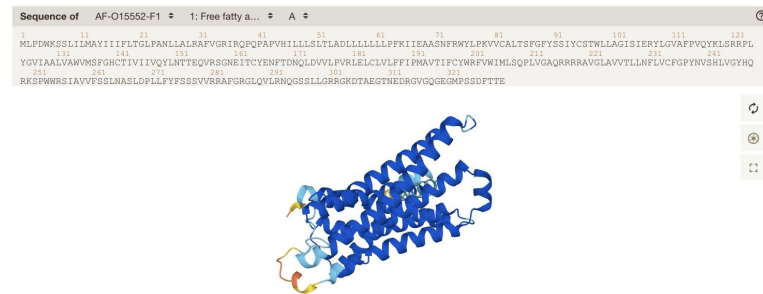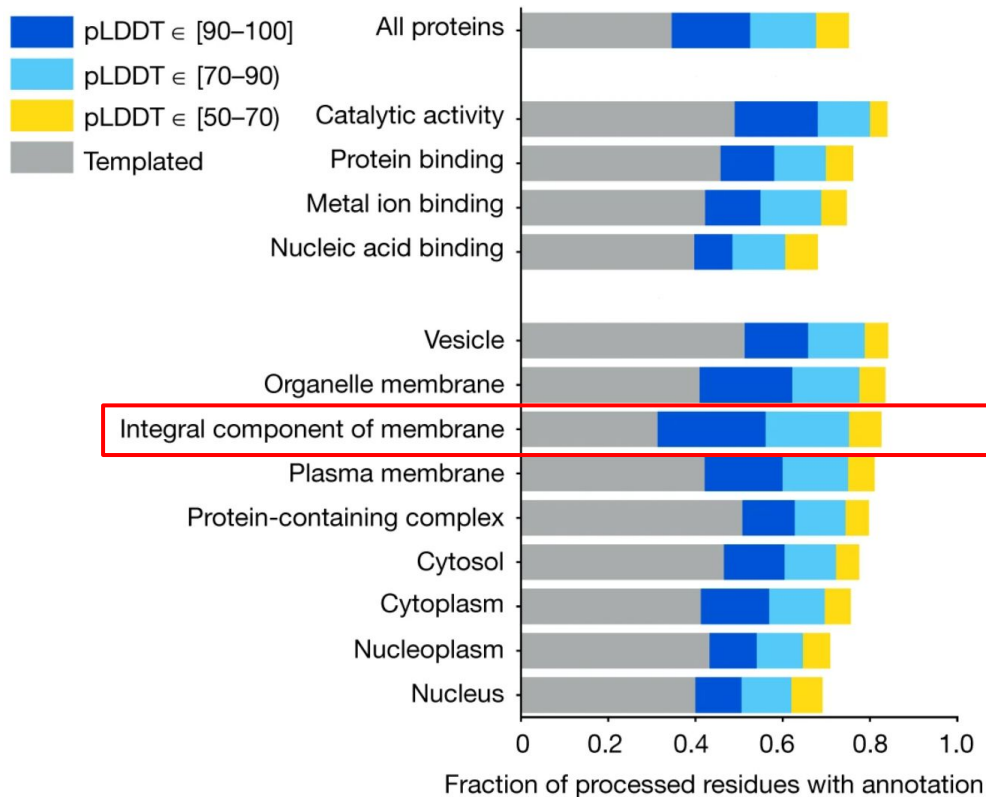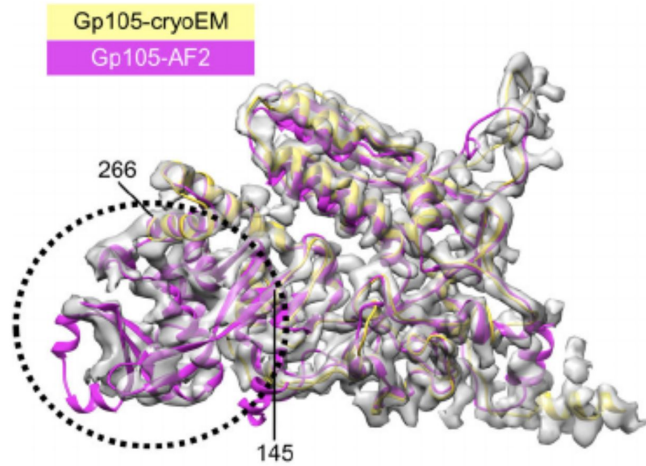- Plans to expand to **Uniref90** (~135M structures)

# Increase in coverage of the human proteome

Even when accounting for template–modelling, AlphaFold greatly expands the high–accuracy structural coverage of the human proteome
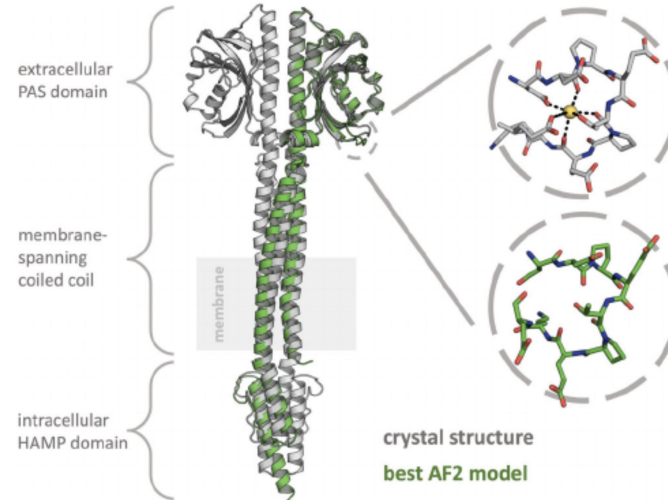
# AlphaFold as an aid to experimental structure determination



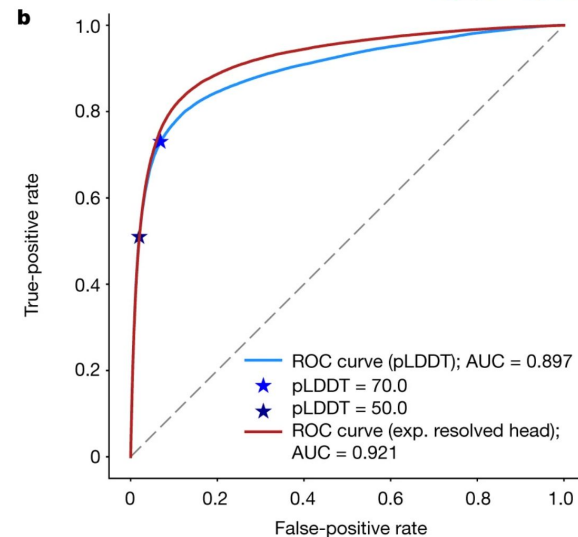**AlphaFold models of AR9 nvRNAP proteins fit the cryo–EM density nearly perfectly**

**Crystal structure of dimeric Af1503**

# AlphaFold confidence and disorder

- In IDRs, there is no fixed structure for AlphaFold to identify
- So AlphaFold produces a boring ribbon *and* reports very low confidence (pLDDT)
- This make very low confidence a strong signal of disorder or chains that are unstructured in isolation



**a**

Fraction of residues

Legend:
- pLDDT > 90
- pLDDT > 70
- pLDDT > 50
- pLDDT ≤ 50

X-axis: PDB resolved, PDB unresolved, Human

**b**

True-positive rate vs False-positive rate

- ROC curve (pLDDT); AUC = 0.897
- ★ pLDDT = 70.0
- ★ pLDDT = 50.0
- ROC curve (exp. resolved head); AUC = 0.921

AlphaFold accuracy on CAID disorder benchmark

**Kamil Górecki** @kamil_gorecki85

You really appreciate AlphaFold when you run it on a protein that for a year refused to get expressed and purified...

10:33 pm · 20 Jul 2021 · Twitter for iPhone

---

**SmithLabUMBC** @SmithLabUMBC

Aaaaaaaand using an AlphaFold model we just phased some very important X-ray data that we previously couldn't phase using other MR approaches and even SAD methods!!!!!

3:46 PM · Jul 23, 2021 · Twitter Web App

---

**Tristan Croll** @CrollTristan · Jul 23

Upshot: while AlphaFold clearly isn't a *replacement* for experimental structures by any stretch, it's already very clear that it's going to make the task of *building* experimental structures both much easier and much less error prone. Welcome to the future! (fin)
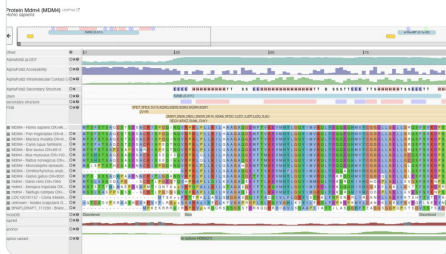
3          35          287

---

**Bálint Mészáros** @_BalintMeszaros

ProViz by @DaveyLab now has various AlphaFold-based disordered predictions and secondary structures

Short Linear Motif team @DaveyLab · Aug 3
@_BalintMeszaros and I have been staring at the @DeepMind @emblebi structure database non-stop. To see them in context we've updated ProViz to visualise AlphaFold2 data mapped to multiple sequence alignments and data from @uniprot @rcsbPDB @PfamDB tinyurl.com/ProVizAF.
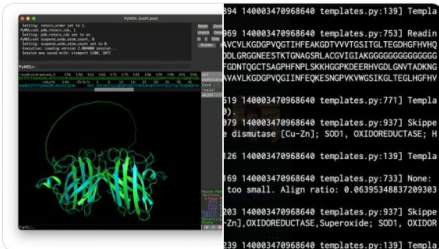
Show this thread

12:49 PM · Aug 3, 2021 · Twitter Web App

---

**Yoshitaka Moriwaki** @Ag_smith

Prediction from AlphaFold2, SOD1 (Superoxide dismutase), max_template_date=1979-07-19 (No template).

Nevertheless AF2 could successfully predict the homodimer form if we input two SOD1 sequences with a polyG linker.
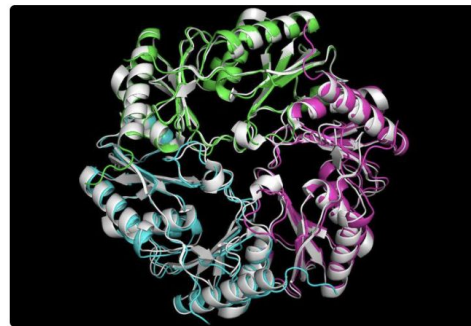
11:39 pm · 19 Jul 2021 · Twitter Web App

---

**Sergey Ovchinnikov** @sokrypton

Homooligomeric prediction in #alphafold works a little too good. So far worked on nearly every case we (me & @minkbaek) tried. Going beyond dimers! Seems @DeepMind accidentally "solved" the homooligomeric prediction problem (w/ MSA input) 😂 Give it a try:
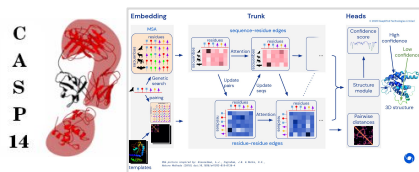https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/AlphaFold2.ipynb
https://pbs.twimg.com/media/E62EYFSXEAIIvCa.jpg

Twitter · Jul 21st (159 kB) ▾
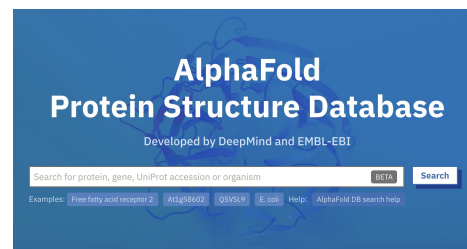
# AlphaFold timeline



### CASP14 conference

### AlphaFold paper & code release

### AlphaFold DB & Proteome paper
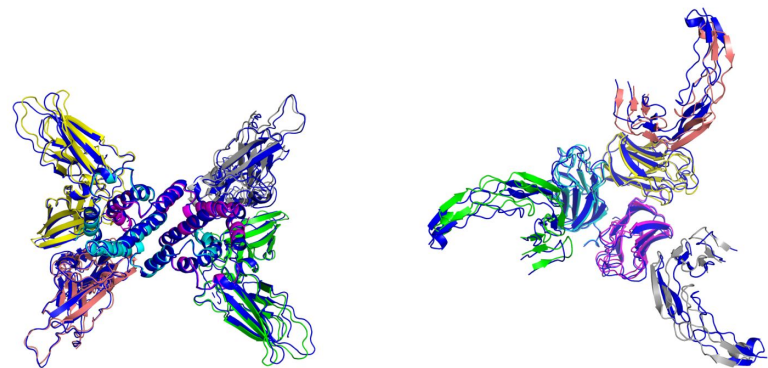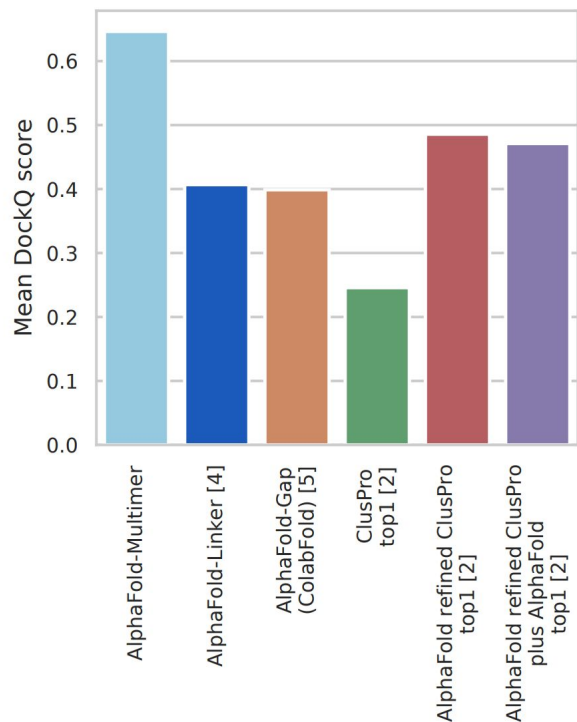
November 2020

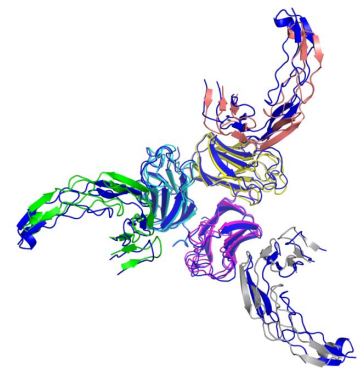July 2021

DeepMind

# AlphaFold–Multimer

# Training AlphaFold to predict protein complexes (AlphaFold-Multimer)

Adapting the inputs, loss function, and training of AlphaFold to handle multimers and then training the model from scratch



(a) A2B2C2 heteromer
TM-score = 98.0, $N_{res}$ = 1,246, PDB ID = 6E3K

(b) A3B3 heteromer
TM-score = 89.3, $N_{res}$ = 795, PDB ID = 7KHD

(c) Protein-peptide complex
TM-score = 96.0, DockQ = 0.948, $N_{res}$ = 385, PDB ID = 6JMT

(d) A2B2 heteromer
TM-score = 98.3, $N_{res}$ = 716, PDB ID = 6IWD

Richard Evans et al. 2021 - https://doi.org/10.1101/2021.10.04.463034

# Organizing

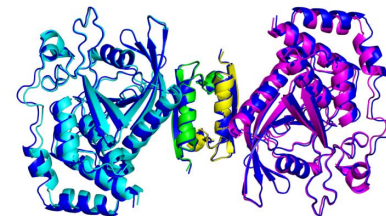- Make it easy to measure your performance – one metric, one leaderboard

- Incrementally improve this metric

  - +0.5% per week adds up – enough quantitative change adds up to qualitative change!

- How to create research velocity:

  - Enable fast iteration – the more ideas you can test, the faster you progress

  - Allow people to build on each other's work – you should always be improving SOTA

  - Test your code – avoid errors that damage progress

# DeepMind work in science

- Protein structure prediction (AlphaFold)

- Quantum chemistry (QMC, DFT)

- Genomics

- Weather prediction

- Fusion reactor control

- Lattice QCD

- Glassy dynamics

- Mathematical discovery

DeepMind

# Acknowledgements

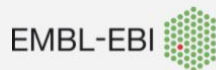# Thank you to everyone who made AlphaFold possible!

Agata Laydon
Alex Bateman
Alex Bridgland
Alexander Pritzel
Andrew Cowie
Andrew J. Ballard
Andrew W. Senior
Anna Potapenko
Augustin Žídek
Bernardino Romera–Paredes

Clemens Meyer
David Reiman
David Silver
Demis Hassabis
Ellen Clancy
Ewan Birney
Gerard J. Kleywegt
John Jumper
Jonas Adler
Kathryn Tunyasuvunakool

Koray Kavukcuoglu
Martin Steinegger
Michael Figurnov
Michal Zielinski
Michalina Pacholska
Olaf Ronneberger
Oriol Vinyals
Pushmeet Kohli
Richard Evans
Rishub Jain

Russ Bates
Sameer Velankar
Sebastian Bodenstein
Simon A. A. Kohl
Stanislav Nikolov
Stig Petersen
Tamas Berghammer
Tim Green
Trevor Back
Zachary Wu

**The wider team at DeepMind and EMBL–EBI**

**The CASP community**     **The experimental biology community**

DeepMind

Q&A