DeepMind

## Highly accurate protein structure prediction with AlphaFold

Michael Figurnov Staff Research Scientist, DeepMind

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

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

## Agenda

- Introduction to protein structure
- How AlphaFold works
- AlphaFold and the biological community
- What's next?
- Lattice QCD work at DeepMind
- Q&A



DeepMind

# Introduction to protein structure



## What is a protein?

- 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





## **Protein structure: terminology**

∇\$8∞0\$# ┃▷⋈□▷∷◇ ○Ѧ屳○₩๏\$

- Vocabulary: 20 common amino acids / residues
  - Small organic molecules with common groups
    - + a side chain specific to each amino acid



Sequence: chain of 100s-1000s amino acids

- Amino acids form peptide bonds and build up a protein chain
- DNA sequences directly encode the amino acid sequence



Structure: atom coordinates in 3D space (300 - 50,000 atoms)

Unique 3D structure comes from physical interactions of amino acids.



## **Protein structure**



## Why predict structures?

Predicting a protein's structure from its amino acid sequence has been a **grand challenge in biology** for the past 50 years.

- → Experimental structure determination takes months to years.
- → ~200,000 protein structures experimentally determined so far.
- → Structure prediction can provide actionable information faster.



## **Reading DNA has become cheap**

#### Cost of sequencing a full human genome



The cost of sequencing the DNA of the complete human genome, measured in US\$. This data is not adjusted for inflation.



200,000x decrease in 20 years

# AlphaFold predicts highly accurate protein structures from amino acid sequences

SIFSYITESTGTPSNATYT YVIERWDPETSGILNPCYG WPVCYVTVNHKHTVNGTGG NPAFQIARIEKLRTLAEVR DVVLKNRSFPIEGQTTHRG PSLNSNQECVGLFYQPNSS GISPRGKLLPGSLCGIAPP PVHHHHH





**T1049 / 6y4f** 93.5 GDT (adhesin tip)





External independent benchmarks are critical Run every 2 years since 1994 – gold standard Blind prediction assessment

CASP

Critical Assessment of protein Structure Prediction



## AlphaFold at CASP14







**T1037 / 6vr4** 90.7 GDT, RNA polymerase domain



## Making AlphaFold available

Article

Highly accurate protein structure prediction with AlphaFold

https://doi.org/10.1038/s41586-021-03819-2 John Jumper<sup>14</sup>, Richard Evans<sup>14</sup>, Alexander Pritzel<sup>14</sup>, Tim Green<sup>14</sup>, Michael Figurnov<sup>14</sup>, Olaf Ronneberger<sup>14</sup>, Kathryn Tunyasuvunakool<sup>14</sup>, Russ Bates<sup>14</sup>, Augustin Židek<sup>14</sup> Received: 11 May 2021 Anna Potapenko<sup>14</sup>, Alex Bridgland<sup>14</sup>, Clemens Meyer<sup>14</sup>, Simon A. A. Kohl<sup>14</sup>, Accepted: 12 July 2021 Andrew J. Ballard<sup>14</sup>, Andrew Cowie<sup>14</sup>, Bernardino Romera-Paredes<sup>14</sup>, Stanislav Nikolov Rishub Jain<sup>14</sup>, Jonas Adler<sup>1</sup>, Trevor Back<sup>1</sup>, Stig Petersen<sup>1</sup>, David Reiman<sup>1</sup>, Ellen Clancy<sup>1</sup>, Published online: 15 July 2021 Michal Zielinski<sup>1</sup>, Martin Steinegger<sup>2,3</sup>, Michalina Pacholska<sup>1</sup>, Tamas Berghammer<sup>1</sup>, Sebastian Bodenstein<sup>1</sup>, David Silver<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuogl Open access Pushmeet Kohli<sup>1</sup> & Demis Hassabis<sup>141</sup> Check for updates

> Proteins are essential to life, and understanding their structure can facilitate a mechanistic understanding of their function. Through an enormous experimental effort1-4, the structures of around 100,000 unique proteins have been determined5, but this represents a small fraction of the billions of known protein sequences<sup>67</sup>. Structural coverage is bottlenecked by the months to years of painstaking effort required to determine a single protein structure. Accurate computational approaches are needed to address this gap and to enable large-scale structural bioinformatics. Predicting the three-dimensional structure that a protein will adopt based solely on its amino acid sequence-the structure prediction component of the 'protein folding problem'8-has been an important open research problem for more than 50 years? Despite recent progress<sup>10-14</sup>, existing methods fall far short of atomic accuracy, especially when no homologous structure is available. Here we provide the first computational method that can regularly predict protein structures with atomic accuracy even in cases in which no similar structure is known. We validated an entirely redesigned version of our neural network-based model, AlphaFold, in the challenging 14th Critical Assessment of protein Structure Prediction (CASP14)15, demonstrating accuracy competitive with experimental structures in a majority of cases and greatly outperforming other methods. Underpinning the latest version of AlphaFold is a novel machine learning approach that incorporates physical and biological knowledge about protein structure. leveraging multi-sequence alignments, into the design of the deep learning algorithm.

The development of computational methods to predict the steady growth of experimental protein structures deposited in three-dimensional (3D) protein structures from the protein sequence the Protein Data Bank (PDB)<sup>5</sup>, the explosion of genomic sequencing has proceeded along two complementary paths that focus on either the and the rapid development of deep learning techniques to interpre physical interactions or the evolutionary history. The physical interactions correlations. Despite these advances, contemporary physical tion programme heavily integrates our understanding of molecular and evolutionary history-based approaches produce predictions that driving forces into either thermodynamic or kinetic simulation of protein physics<sup>16</sup> or statistical approximations thereof<sup>17</sup>. Although theoreti a close homologue has not been solved experimentally and this has cally very appealing, this approach has proved highly challenging for limited their utility for many biological applications even moderate-sized proteins due to the computational intractability of molecular simulation, the context dependence of protein stability approach capable of predicting protein structures to near experimental and the difficulty of producing sufficiently accurate models of protein accuracy in a majority of cases. The neural network AlphaFold that we physics. The evolutionary programme has provided an alternative in developed was entered into the CASP14 assessment (May-July 2020; recent years in which the constraints on protein structure are derived entered under the team name 'AlphaFold2' and a completely different from bioinformatics analysis of the evolutionary history of proteins. model from our CASP13 AlphaFold system<sup>30</sup>). The CASP assessment is homology to solved structures<sup>18,30</sup> and nairwise evolutionary correla-

In this study, we develop the first, to our knowledge, computational tions<sup>30-34</sup>. This bioinformatics approach has benefited greatly from been deposited in the PDB or publicly disclosed so that it is a blind test

DeepMind, London, UK. "School of Biological Sciences, Seoul National University, Seoul, South Korea, "Artificial Intelligence Institute, Seoul National University, Seoul, South Korea, "These authors contributed equally, john tumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Flournov, Olaf Ronneberger, Kathyn Tumvayuyunakool, Buss Rates, Augustin Židek, Anna tapenio, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowle, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Demis Hassabi <sup>16</sup>e-mail: lumpenRdeepmind.com/dbcontactRdeepmind.com

Nature | Vol 596 | 26 August 2021 | 583

#### deepmind / alphafold Public

公

ဗ္ Fork 1.7k Star 9.6k

BETA

1

Search



Search for protein, gene, UniProt accession or organism

Examples: Free fatty acid receptor 2 At1g58602 Q5VSL9 E. coli Help: AlphaFold DB search help

#### >200 million protein structures



DeepMind

# How AlphaFold works

## Why not physical simulations?





## Why not physical simulations? Reason 1: timescales

- Molecular dynamics (MD):
  - Langevin dynamics on a potential
  - Folding takes 10<sup>12</sup>–10<sup>15</sup> sequential steps too slow





SARS-COV-2 spike, Anton simulation

|--|--|--|--|--|--|

## Why not physical simulations? Reason 2: underspecified context





## **Determining Structure from Evolution - Intuition**





## **Determining Structure from Evolution - Intuition**





Coevolution cartoon by Sergey Ovchinnikov (https://jgi.doe.gov/seeking-structure-metagenome -sequences/cartoon-coevolution-sergey-o/) Deep learning provides building blocks for approximating arbitrary functions



These blocks are very generic; they don't take advantage of our scientific understanding of proteins and protein evolution



## **Design principles**

End-to-end

- Network goes all the way from inputs to structure
- Can learn about and optimise the whole process

#### Inductive biases

- Design reflects our knowledge of physics / geometry
- Emphasis on pairs, not a sequence of residues
- Output should be self-consistent





## AlphaFold: a deep learning model

#### Model

- Deep: 192 blocks (48 blocks x 4 cycles)
- 93 million parameters
  - Reasonably small; GPT-3 is 175 *billion* parameters
- Seconds to hours of runtime, depending on protein length

#### **Training Data**

• 170,000 3D structures



• Genetic databases (sequences without structure)



## **Model overview**



Incorporates evolutionary information

Determines relationship between residue pairs End-to-end: outputs structure directly











Attention is augmented by the network's belief about residue pairs





Outer product allows generalized correlation similar to co-evolution





Residue pair interactions inspired by geometric interactions

## **Triangular attention for residue pairs**

- → Matrix view:
  - Each residue attends to all other residues within a row



Corresponding edges in a graph





## **Triangular attention for residue pairs**

- → Matrix view:
  - Each residue attends to all other residues within a row



#### → Graph view:

- Nodes = residues, edges = pairs of residues; represents residue spatial relations, e.g. distances in 3D.
- Transitivity / triangle inequality: update for query ij should depend on ik and jk (for residues i, j, k).

$$a_{ijk}^{h} = \operatorname{softmax}_{k} \left( \frac{1}{\sqrt{c}} \mathbf{q}_{ij}^{h^{\top}} \mathbf{k}_{ik}^{h} + b_{jk}^{h} \right)$$

This is missing in traditional self-attention: we add a bias term to bring the third edge in.



## Network



6

## **Structure module**



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



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



## **Geometric transformer architecture** updates the rigid bodies / backbone

• Also builds the side chains from torsion angles



## **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**





6VR4: Drobysheva, A.V., et al. Structure and function of virion RNA polymerase of a crAss-like phage. Nature (2020). (CASP14 target T1044) 1

## **Model interpretability - Role of depth**



DeepMind

# AlphaFold and the biology community

#### 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



ProViz by @DaveyLab now has various AlphaFoldbased 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



Pedro Beltrao @pedrobeltrao

We joined a large community effort to assess diverse applications of AlphaFold 2 in the context of novel structural elements; missense variants; function and ligand binding sites; modelling of interactions and experimental structural data. Some highlights below:

#### bioRxiv @biorxivpreprint · 12h

A structural biology community assessment of AlphaFold 2 applications biorxiv.org/cgi/content/sh... #bioRxiv

7:35 AM · Sep 27, 2021 · Twitter Web App



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



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)



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) is Give it a try: https://colab.research.google.com/github/sokrypton/Col abFold/blob/main/AlphaFold2.ipynb https://pbs.twimg.com/media/E62EYFSXEAIIvCa.jpg Twitter |Jul 21st (159 kB) •



## **Communicating confidence is key**

- AlphaFold outputs two confidence metrics: **pLDDT** and **PAE**
- Important for sharing predictions responsibly
- Supports new use cases:
  - Software integrations
  - Domain segmentation
  - Detecting interactions
  - Disorder prediction
  - Ranking structural models . . .



## **Studying molecular machines**



Mosalaganti, S. et al. Al-based structure prediction empowers integrative structural analysis of human nuclear pores. Science (2022)



DeepMind

# What's next?

## What's next? Three major challenges

#### Proteins in the cellular context

- AlphaFold-Multimer: multiple protein chains
- Need to also consider non-protein components: DNA, RNA, ligands, water, ions...





## What's next? Three major challenges

#### Proteins in the cellular context

- AlphaFold–Multimer: multiple protein chains
- Need to also consider non-protein components: DNA, RNA, ligands, water, ions...

#### **Protein dynamics**

#### Effect of mutations on proteins









DESRES-ANTON-11021571 https://www.deshawresearch.com/downloads/download\_trajectory\_sarscov2.cgi/ RNA polymerase, PDB-101 https://pdb101.rcsb.org/motm/40

## Thank you to everyone who made AlphaFold possible!

#### AlphaFold 2 Methods

John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger (Seoul NU), Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, Demis Hassabis

#### Human Proteome

Kathryn Tunyasuvunakool, Jonas Adler, Zachary Wu, Tim Green, Michal Zielinski, Augustin Žídek, Alex Bridgland, Andrew Cowie, Clemens Meyer, Agata Laydon, Sameer Velankar (EBI), Gerard J. Kleywegt (EBI), Alex Bateman (EBI), Richard Evans, Alexander Pritzel, Michael Figurnov, Olaf Ronneberger, Russ Bates, Simon Kohl, Anna Potapenko, Andrew J. Ballard, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Ellen Clancy, David Reiman, Stig Petersen, Andrew W. Senior, Koray Kavukcuoglu, Ewan Birney (EBI), Pushmeet Kohli, John Jumper, Demis Hassabis



DeepMind

# Lattice QCD work at DeepMind

## Enabling Lattice QCD simulations to a wide community and accelerate progress in our understanding of matter

#### The problem; scientific viewpoint:

- Simulate quantum fluctuations of gauge fields (gluons, photons) and fermionic fields (quarks, electrons) in QFT using generative models
- Use the model's samples to compute physical observables (such as particle masses)

#### The problem; machine learning viewpoint:

- Not a data fitting problem!
- Inference problem:
  - Learn a target density in a 4D lattice of complex-valued matrices and vectors
  - Minimize a divergence using model sampler (e.g. reverse Kullback-Leibler)
- Requires enormous amount of memory due to lattice size



 $p(U) \propto e^{-\beta S[U]}$ 



C DeepMind

6

Sampling using SU(N) gauge equivariant flows, Phys. Rev. D, Boyda, Kanwar, Racaniere, Rezende, Albergo, Cranmer, Hackett, Shanahan
Gauge-equivariant flow models for sampling in lattice field theories with pseudofermions, Phys. Rev. D, Abbott, Albergo, Boyda,

Cranmer, Hackett, Kanwar, Racanière, Rezende, Romero-López, Shanahan, Tian, Urban

DeepMind

# Thank you!

Q&A