

# Highly accurate protein structure prediction with AlphaFold

feature: dataset search + MSA + self dot pro

<https://doi.org/10.1038/s41586-021-03819-2>



Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

Open access

 Check for updates

John Jumper<sup>1,4</sup>, Richard Evans<sup>1,4</sup>, Alexander Pritzel<sup>1,4</sup>, Tim Green<sup>1,4</sup>, Michael Figurnov<sup>1,4</sup>, Olaf Ronneberger<sup>1,4</sup>, Kathryn Tunyasuvunakool<sup>1,4</sup>, Russ Bates<sup>1,4</sup>, Augustin Židek<sup>1,4</sup>, Anna Potapenko<sup>1,4</sup>, Alex Bridgland<sup>1,4</sup>, Clemens Meyer<sup>1,4</sup>, Simon A. A. Kohl<sup>1,4</sup>, Andrew J. Ballard<sup>1,4</sup>, Andrew Cowie<sup>1,4</sup>, Bernardino Romera-Paredes<sup>1,4</sup>, Stanislav Nikolov<sup>1,4</sup>, Rishub Jain<sup>1,4</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>, Michal Zielinski<sup>1</sup>, Martin Steinegger<sup>2,3</sup>, Michalina Pacholska<sup>1</sup>, Tamas Berghammer<sup>1</sup>, Sebastian Bodenstern<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>1,4</sup>

Proteins are essential to life, and understanding their structure can facilitate a mechanistic understanding of their function. Through an enormous experimental effort<sup>1–4</sup>, the structures of around 100,000 unique proteins have been determined<sup>5</sup>, but this represents a small fraction of the billions of known protein sequences<sup>6,7</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’<sup>8</sup>—has been an important open research problem for more than 50 years<sup>9</sup>. 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)<sup>15</sup>, 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.

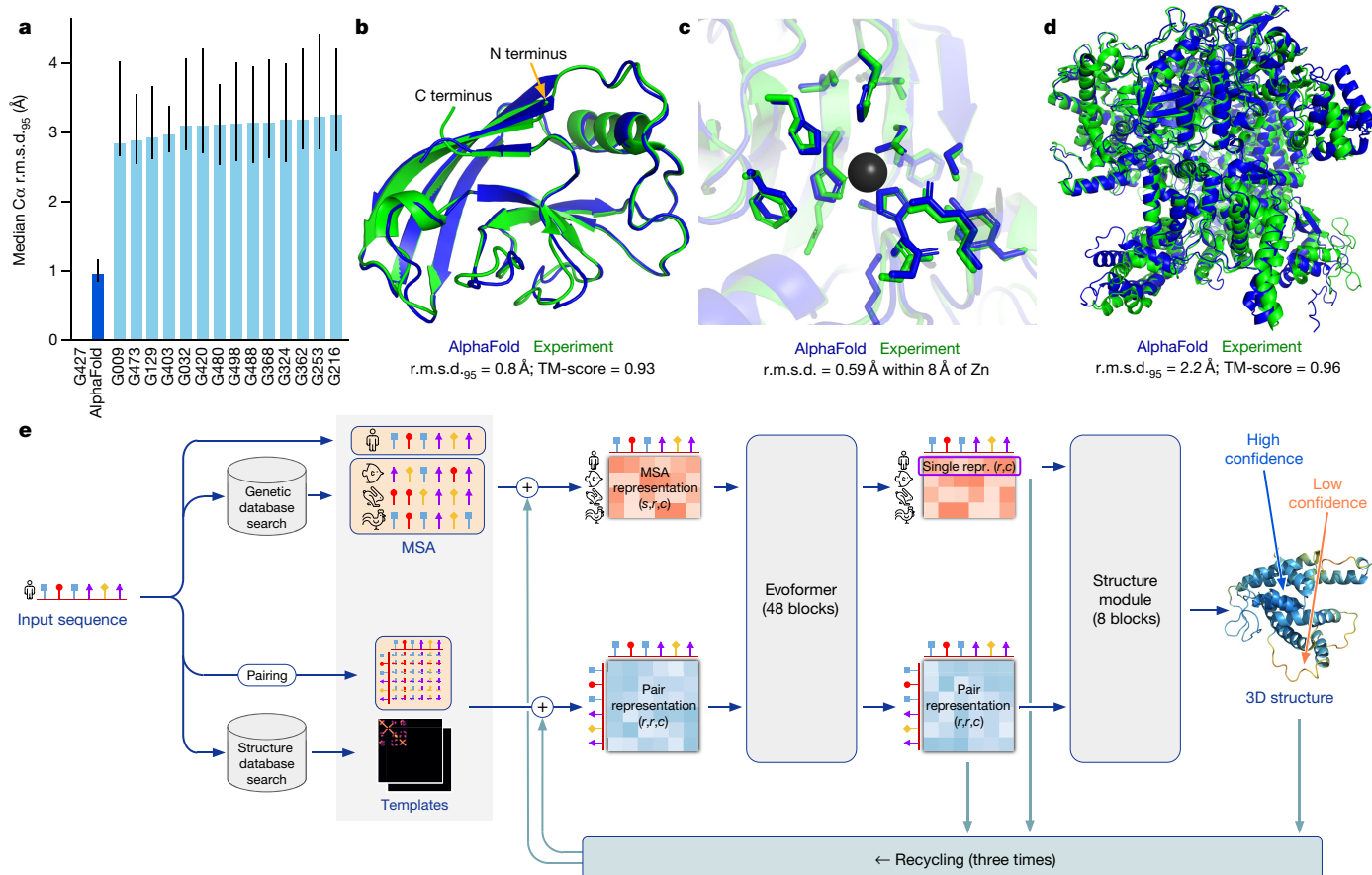
The development of computational methods to predict three-dimensional (3D) protein structures from the protein sequence has proceeded along two complementary paths that focus on either the physical interactions or the evolutionary history. The physical interaction programme heavily integrates our understanding of molecular driving forces into either thermodynamic or kinetic simulation of protein physics<sup>16</sup> or statistical approximations thereof<sup>17</sup>. Although theoretically very appealing, this approach has proved highly challenging for even moderate-sized proteins due to the computational intractability of molecular simulation, the context dependence of protein stability and the difficulty of producing sufficiently accurate models of protein physics. The evolutionary programme has provided an alternative in recent years, in which the constraints on protein structure are derived from bioinformatics analysis of the evolutionary history of proteins, homology to solved structures<sup>18,19</sup> and pairwise evolutionary correlations<sup>20–24</sup>. This bioinformatics approach has benefited greatly from

the steady growth of experimental protein structures deposited in the Protein Data Bank (PDB)<sup>5</sup>, the explosion of genomic sequencing and the rapid development of deep learning techniques to interpret these correlations. Despite these advances, contemporary physical and evolutionary-history-based approaches produce predictions that are far short of experimental accuracy in the majority of cases in which a close homologue has not been solved experimentally and this has limited their utility for many biological applications.

In this study, we develop the first, to our knowledge, computational approach capable of predicting protein structures to near experimental accuracy in a majority of cases. The neural network AlphaFold that we developed was entered into the CASP14 assessment (May–July 2020; entered under the team name ‘AlphaFold2’ and a completely different model from our CASP13 AlphaFold system<sup>10</sup>). The CASP assessment is carried out biennially using recently solved structures that have not been deposited in the PDB or publicly disclosed so that it is a blind test

<sup>1</sup>DeepMind, London, UK. <sup>2</sup>School of Biological Sciences, Seoul National University, Seoul, South Korea. <sup>3</sup>Artificial Intelligence Institute, Seoul National University, Seoul, South Korea. <sup>4</sup>These authors contributed equally: John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Demis Hassabis.

<sup>✉</sup>e-mail: jumper@deepmind.com; dhcontact@deepmind.com



**Fig. 1 | AlphaFold produces highly accurate structures.** **a**, The performance of AlphaFold on the CASP14 dataset ( $n = 87$  protein domains) relative to the top-15 entries (out of 146 entries), group numbers correspond to the numbers assigned to entrants by CASP. Data are median and the 95% confidence interval of the median, estimated from 10,000 bootstrap samples. **b**, Our prediction of CASP14 target T1049 (PDB 6Y4F, blue) compared with the true (experimental) structure (green). Four residues in the C terminus of the crystal structure are B-factor outliers and are not depicted. **c**, CASP14 target T1056 (PDB 6YJ1).

An example of a well-predicted zinc-binding site (AlphaFold has accurate side chains even though it does not explicitly predict the zinc ion). **d**, CASP target T1044 (PDB 6VR4)—a 2,180-residue single chain—was predicted with correct domain packing (the prediction was made after CASP using AlphaFold without intervention). **e**, Model architecture. Arrows show the information flow among the various components described in this paper. Array shapes are shown in parentheses with  $s$ , number of sequences ( $N_{\text{seq}}$  in the main text);  $r$ , number of residues ( $N_{\text{res}}$  in the main text);  $c$ , number of channels.

for the participating methods, and has long served as the gold-standard assessment for the accuracy of structure prediction<sup>25,26</sup>.

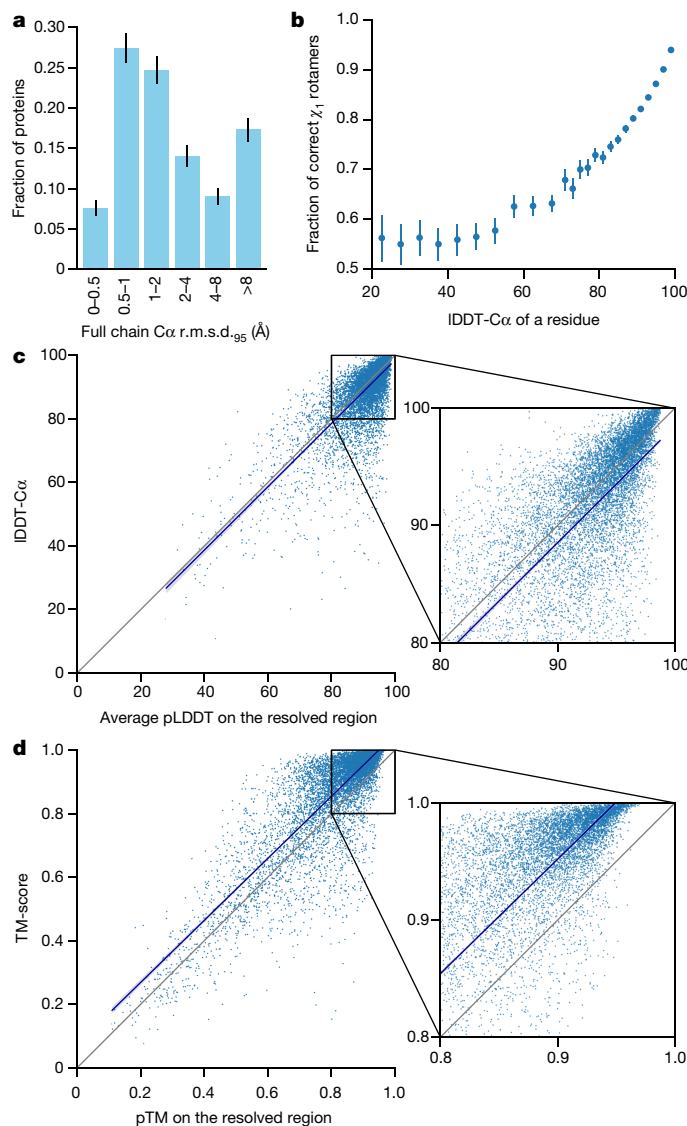
In CASP14, AlphaFold structures were vastly more accurate than competing methods. AlphaFold structures had a median backbone accuracy of 0.96 Å r.m.s.d.<sub>95</sub> (C $\alpha$  root-mean-square deviation at 95% residue coverage) (95% confidence interval = 0.85–1.16 Å) whereas the next best performing method had a median backbone accuracy of 2.8 Å r.m.s.d.<sub>95</sub> (95% confidence interval = 2.7–4.0 Å) (measured on CASP domains; see Fig. 1a for backbone accuracy and Supplementary Fig. 14 for all-atom accuracy). As a comparison point for this accuracy, the width of a carbon atom is approximately 1.4 Å. In addition to very accurate domain structures (Fig. 1b), AlphaFold is able to produce highly accurate side chains (Fig. 1c) when the backbone is highly accurate and considerably improves over template-based methods even when strong templates are available. The all-atom accuracy of AlphaFold was 1.5 Å r.m.s.d.<sub>95</sub> (95% confidence interval = 1.2–1.6 Å) compared with the 3.5 Å r.m.s.d.<sub>95</sub> (95% confidence interval = 3.1–4.2 Å) of the best alternative method. Our methods are scalable to very long proteins with accurate domains and domain-packing (see Fig. 1d for the prediction of a 2,180-residue protein with no structural homologues). Finally, the model is able to provide precise, per-residue estimates of its reliability that should enable the confident use of these predictions.

We demonstrate in Fig. 2a that the high accuracy that AlphaFold demonstrated in CASP14 extends to a large sample of recently released PDB

structures; in this dataset, all structures were deposited in the PDB after our training data cut-off and are analysed as full chains (see Methods, Supplementary Fig. 15 and Supplementary Table 6 for more details). Furthermore, we observe high side-chain accuracy when the backbone prediction is accurate (Fig. 2b) and we show that our confidence measure, the predicted local-distance difference test (pLDDT), reliably predicts the C $\alpha$  local-distance difference test (IDDT-C $\alpha$ ) accuracy of the corresponding prediction (Fig. 2c). We also find that the global superposition metric template modelling score (TM-score)<sup>27</sup> can be accurately estimated (Fig. 2d). Overall, these analyses validate that the high accuracy and reliability of AlphaFold on CASP14 proteins also transfers to an uncurated collection of recent PDB submissions, as would be expected (see Supplementary Methods 1.15 and Supplementary Fig. 11 for confirmation that this high accuracy extends to new folds).

## The AlphaFold network

AlphaFold greatly improves the accuracy of structure prediction by incorporating novel neural network architectures and training procedures based on the evolutionary, physical and geometric constraints of protein structures. In particular, we demonstrate a new architecture to jointly embed multiple sequence alignments (MSAs) and pairwise features, a new output representation and associated loss that enable accurate end-to-end structure prediction, a new equivariant attention



**Fig. 2 | Accuracy of AlphaFold on recent PDB structures.** The analysed structures are newer than any structure in the training set. Further filtering is applied to reduce redundancy (see Methods). **a**, Histogram of backbone r.m.s.d. for full chains (C $\alpha$  r.m.s.d. at 95% coverage). Error bars are 95% confidence intervals (Poisson). This dataset excludes proteins with a template (identified by *hmmsearch*) from the training set with more than 40% sequence identity covering more than 1% of the chain ( $n = 3,144$  protein chains). The overall median is 1.46 Å (95% confidence interval = 1.40–1.56 Å). Note that this measure will be highly sensitive to domain packing and domain accuracy; a high r.m.s.d. is expected for some chains with uncertain packing or packing errors. **b**, Correlation between backbone accuracy and side-chain accuracy. Filtered to structures with any observed side chains and resolution better than 2.5 Å ( $n = 5,317$  protein chains); side chains were further filtered to  $B$ -factor  $< 30$  Å<sup>2</sup>. A rotamer is classified as correct if the predicted torsion angle is within 40°. Each point aggregates a range of IDDT-C $\alpha$ , with a bin size of 2 units above 70 IDDT-C $\alpha$  and 5 units otherwise. Points correspond to the mean accuracy; error bars are 95% confidence intervals (Student  $t$ -test) of the mean on a per-residue basis. **c**, Confidence score compared to the true accuracy on chains. Least-squares linear fit IDDT-C $\alpha = 0.997 \times \text{pLDDT} - 1.17$  (Pearson's  $r = 0.76$ ).  $n = 10,795$  protein chains. The shaded region of the linear fit represents a 95% confidence interval estimated from 10,000 bootstrap samples. In the companion paper<sup>39</sup>, additional quantification of the reliability of pLDDT as a confidence measure is provided. **d**, Correlation between pTM and full chain TM-score. Least-squares linear fit TM-score =  $0.98 \times \text{pTM} + 0.07$  (Pearson's  $r = 0.85$ ).  $n = 10,795$  protein chains. The shaded region of the linear fit represents a 95% confidence interval estimated from 10,000 bootstrap samples.

architecture, use of intermediate losses to achieve iterative refinement of predictions, masked MSA loss to jointly train with the structure, learning from unlabelled protein sequences using self-distillation and self-estimates of accuracy.

The AlphaFold network directly predicts the 3D coordinates of all heavy atoms for a given protein using the primary amino acid sequence and aligned sequences of homologues as inputs (Fig. 1e; see Methods for details of inputs including databases, MSA construction and use of templates). A description of the most important ideas and components is provided below. The full network architecture and training procedure are provided in the Supplementary Methods.

The network comprises two main stages. First, the trunk of the network processes the inputs through repeated layers of a novel neural network block that we term Evoformer to produce an  $N_{\text{seq}} \times N_{\text{res}}$  array ( $N_{\text{seq}}$ , number of sequences;  $N_{\text{res}}$ , number of residues) that represents a processed MSA and an  $N_{\text{res}} \times N_{\text{res}}$  array that represents residue pairs. The MSA representation is initialized with the raw MSA (although see Supplementary Methods 1.2.7 for details of handling very deep MSAs). The Evoformer blocks contain a number of attention-based and non-attention-based components. We show evidence in 'Interpreting the neural network' that a concrete structural hypothesis arises early within the Evoformer blocks and is continuously refined. The key innovations in the Evoformer block are new mechanisms to exchange information within the MSA and pair representations that enable direct reasoning about the spatial and evolutionary relationships.

The trunk of the network is followed by the structure module that introduces an explicit 3D structure in the form of a rotation and translation for each residue of the protein (global rigid body frames). These representations are initialized in a trivial state with all rotations set to the identity and all positions set to the origin, but rapidly develop and refine a highly accurate protein structure with precise atomic details. Key innovations in this section of the network include breaking the chain structure to allow simultaneous local refinement of all parts of the structure, a novel equivariant transformer to allow the network to implicitly reason about the unrepresented side-chain atoms and a loss term that places substantial weight on the orientational correctness of the residues. Both within the structure module and throughout the whole network, we reinforce the notion of iterative refinement by repeatedly applying the final loss to outputs and then feeding the outputs recursively into the same modules. The iterative refinement using the whole network (which we term 'recycling' and is related to approaches in computer vision<sup>28,29</sup>) contributes markedly to accuracy with minor extra training time (see Supplementary Methods 1.8 for details).

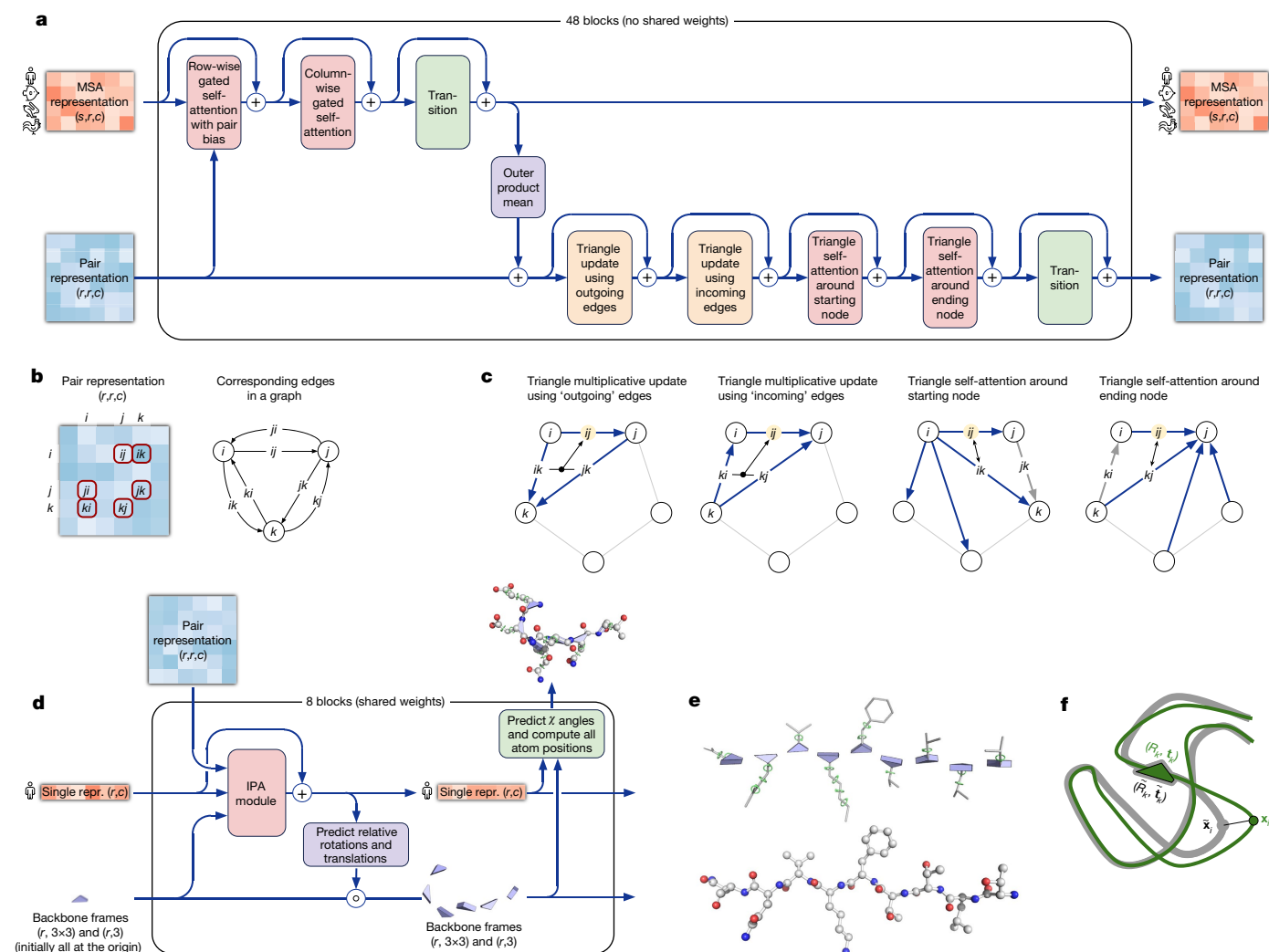
## Evoformer

The key principle of the building block of the network—named Evoformer (Figs. 1e, 3a)—is to view the prediction of protein structures as a graph inference problem in 3D space in which the edges of the graph are defined by residues in proximity. The elements of the pair representation encode information about the relation between the residues (Fig. 3b). The columns of the MSA representation encode the individual residues of the input sequence while the rows represent the sequences in which those residues appear. Within this framework, we define a number of update operations that are applied in each block in which the different update operations are applied in series.

The MSA representation updates the pair representation through an element-wise outer product that is summed over the MSA sequence dimension. In contrast to previous work<sup>30</sup>, this operation is applied within every block rather than once in the network, which enables the continuous communication from the evolving MSA representation to the pair representation.

Within the pair representation, there are two different update patterns. Both are inspired by the necessity of consistency of the pair





**Fig. 3 | Architectural details.** **a**, Evoformer block. Arrows show the information flow. The shape of the arrays is shown in parentheses. **b**, The pair representation interpreted as directed edges in a graph. **c**, Triangle multiplicative update and triangle self-attention. The circles represent residues. Entries in the pair representation are illustrated as directed edges and in each diagram, the edge being updated is  $ij$ . **d**, Structure module including Invariant point attention (IPA)

module. The single representation is a copy of the first row of the MSA representation. **e**, Residue gas: a representation of each residue as one free-floating rigid body for the backbone (blue triangles) and  $\chi$  angles for the side chains (green circles). The corresponding atomic structure is shown below. **f**, Frame aligned point error (FAPE). Green, predicted structure; grey, true structure;  $(R_k, t_k)$ , frames;  $x_i$ , atom positions.

representation—for a pairwise description of amino acids to be representable as a single 3D structure, many constraints must be satisfied including the triangle inequality on distances. On the basis of this intuition, we arrange the update operations on the pair representation in terms of triangles of edges involving three different nodes (Fig. 3c). In particular, we add an extra logit bias to axial attention<sup>31</sup> to include the ‘missing edge’ of the triangle and we define a non-attention update operation ‘triangle multiplicative update’ that uses two edges to update the missing third edge (see Supplementary Methods 1.6.5 for details). The triangle multiplicative update was developed originally as a more symmetric and cheaper replacement for the attention, and networks that use only the attention or multiplicative update are both able to produce high-accuracy structures. However, the combination of the two updates is more accurate.

We also use a variant of axial attention within the MSA representation. During the per-sequence attention in the MSA, we project additional logits from the pair stack to bias the MSA attention. This closes the loop by providing information flow from the pair representation back into the MSA representation, ensuring that the overall Evoformer block is able to fully mix information between the pair and MSA representations and prepare for structure generation within the structure module.

## End-to-end structure prediction

The structure module (Fig. 3d) operates on a concrete 3D backbone structure using the pair representation and the original sequence row (single representation) of the MSA representation from the trunk. The 3D backbone structure is represented as  $N_{\text{res}}$  independent rotations and translations, each with respect to the global frame (residue gas) (Fig. 3e). These rotations and translations—representing the geometry of the N-Cα-C atoms—prioritize the orientation of the protein backbone so that the location of the side chain of each residue is highly constrained within that frame. Conversely, the peptide bond geometry is completely unconstrained and the network is observed to frequently violate the chain constraint during the application of the structure module as breaking this constraint enables the local refinement of all parts of the chain without solving complex loop closure problems. Satisfaction of the peptide bond geometry is encouraged during fine-tuning by a violation loss term. Exact enforcement of peptide bond geometry is only achieved in the post-prediction relaxation of the structure by gradient descent in the Amber<sup>32</sup> force field. Empirically, this final relaxation does not improve the accuracy of the model as measured by the

global distance test (GDT)<sup>33</sup> or IDDT-C $\alpha$ <sup>34</sup> but does remove distracting stereochemical violations without the loss of accuracy.

The residue gas representation is updated iteratively in two stages (Fig. 3d). First, a geometry-aware attention operation that we term ‘invariant point attention’ (IPA) is used to update an  $N_{\text{res}}$  set of neural activations (single representation) without changing the 3D positions, then an equivariant update operation is performed on the residue gas using the updated activations. The IPA augments each of the usual attention queries, keys and values with 3D points that are produced in the local frame of each residue such that the final value is invariant to global rotations and translations (see Methods ‘IPA’ for details). The 3D queries and keys also impose a strong spatial/locality bias on the attention, which is well-suited to the iterative refinement of the protein structure. After each attention operation and element-wise transition block, the module computes an update to the rotation and translation of each backbone frame. The application of these updates within the local frame of each residue makes the overall attention and update block an equivariant operation on the residue gas.

Predictions of side-chain  $\chi$  angles as well as the final, per-residue accuracy of the structure (pLDDT) are computed with small per-residue networks on the final activations at the end of the network. The estimate of the TM-score (pTM) is obtained from a pairwise error prediction that is computed as a linear projection from the final pair representation. The final loss (which we term the frame-aligned point error (FAPE) (Fig. 3f)) compares the predicted atom positions to the true positions under many different alignments. For each alignment, defined by aligning the predicted frame ( $R_i, \mathbf{t}_i$ ) to the corresponding true frame, we compute the distance of all predicted atom positions  $\mathbf{x}_i$  from the true atom positions. The resulting  $N_{\text{frames}} \times N_{\text{atoms}}$  distances are penalized with a clamped  $L^1$  loss. This creates a strong bias for atoms to be correct relative to the local frame of each residue and hence correct with respect to its side-chain interactions, as well as providing the main source of chirality for AlphaFold (Supplementary Methods 1.9.3 and Supplementary Fig. 9).

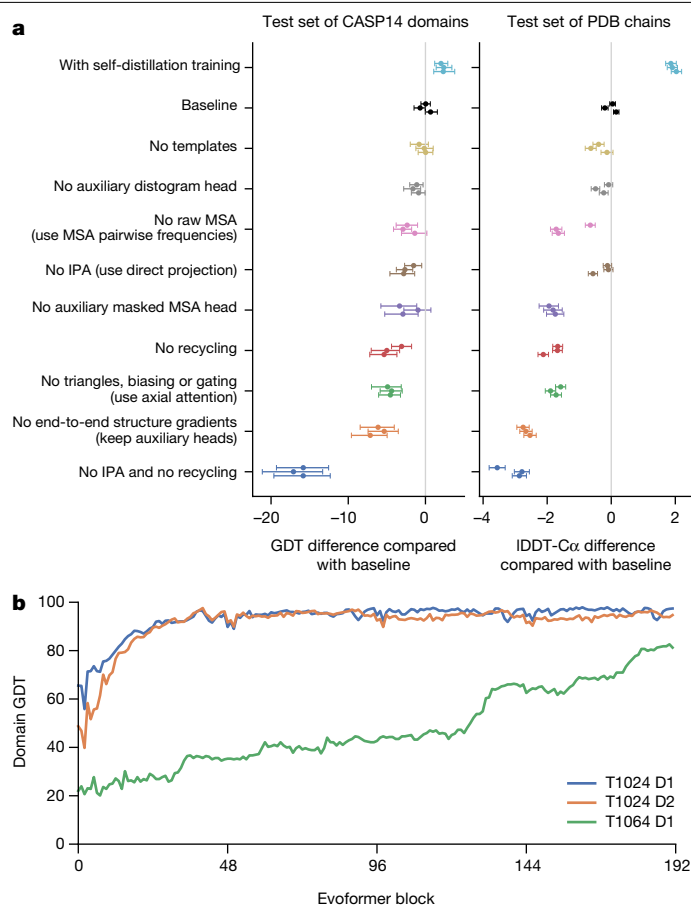
## Training with labelled and unlabelled data

The AlphaFold architecture is able to train to high accuracy using only supervised learning on PDB data, but we are able to enhance accuracy (Fig. 4a) using an approach similar to [noisy student self-distillation](#)<sup>35</sup>. In this procedure, we use a trained network to predict the structure of around 350,000 diverse sequences from Uniclust30<sup>36</sup> and make a new dataset of predicted structures filtered to a high-confidence subset. We then train the same architecture again from scratch using a mixture of PDB data and this new dataset of predicted structures as the training data, in which the various training data augmentations such as cropping and MSA subsampling make it challenging for the network to recapitulate the previously predicted structures. This self-distillation procedure makes effective use of the unlabelled sequence data and considerably improves the accuracy of the resulting network.

Additionally, we randomly [mask out or mutate individual](#) residues within the MSA and have a Bidirectional Encoder Representations from Transformers (BERT)-style<sup>37</sup> objective to predict the masked elements of the MSA sequences. This objective encourages the network to learn to interpret phylogenetic and covariation relationships without hardcoding a particular correlation statistic into the features. The BERT objective is trained jointly with the normal PDB structure loss on the same training examples and is not pre-trained, in contrast to recent independent work<sup>38</sup>.

## Interpreting the neural network

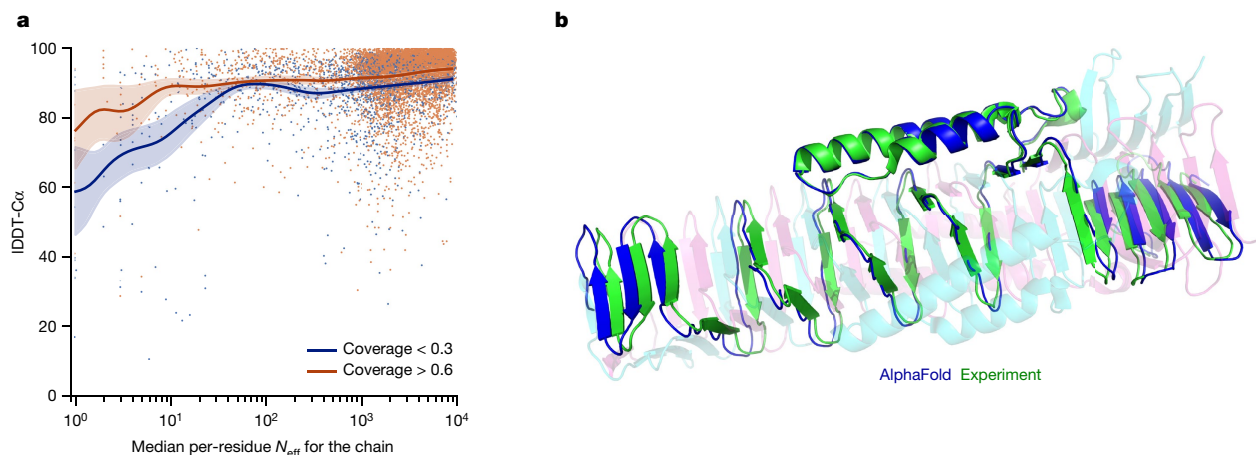
To understand how AlphaFold predicts protein structure, we trained a separate structure module for each of the 48 Evoformer blocks in the network while keeping all parameters of the main network frozen (Supplementary Methods 1.14). Including our recycling stages, this provides a trajectory of 192 intermediate structures—one per full



**Fig. 4 | Interpreting the neural network.** **a**, Ablation results on two target sets: the CASP14 set of domains ( $n = 87$  protein domains) and the PDB test set of chains with template coverage of  $\leq 30\%$  at 30% identity ( $n = 2,261$  protein chains). Domains are scored with GDT and chains are scored with IDDT-C $\alpha$ . The ablations are reported as a difference compared with the average of the three baseline seeds. Means (points) and 95% bootstrap percentile intervals (error bars) are computed using bootstrap estimates of 10,000 samples. **b**, Domain GDT trajectory over 4 recycling iterations and 48 Evoformer blocks on CASP14 targets LmrP (T1024) and Orf8 (T1064) where D1 and D2 refer to the individual domains as defined by the CASP assessment. Both T1024 domains obtain the correct structure early in the network, whereas the structure of T1064 changes multiple times and requires nearly the full depth of the network to reach the final structure. Note, 48 Evoformer blocks comprise one recycling iteration.

Evoformer block—in which each intermediate represents the belief of the network of the most likely structure at that block. The resulting trajectories are surprisingly smooth after the first few blocks, showing that AlphaFold makes constant incremental improvements to the structure until it can no longer improve (see Fig. 4b for a trajectory of accuracy). These trajectories also illustrate the role of network depth. For very challenging proteins such as ORF8 of SARS-CoV-2 (T1064), the network searches and rearranges secondary structure elements for many layers before settling on a good structure. For other proteins such as LmrP (T1024), the network finds the final structure within the first few layers. Structure trajectories of CASP14 targets T1024, T1044, T1064 and T1091 that demonstrate a clear iterative building process for a range of protein sizes and difficulties are shown in Supplementary Videos 1–4. In Supplementary Methods 1.16 and Supplementary Figs. 12, 13, we interpret the attention maps produced by AlphaFold layers.

Figure 4a contains detailed ablations of the components of AlphaFold that demonstrate that a variety of different mechanisms contribute to AlphaFold accuracy. Detailed descriptions of each ablation model, their training details, extended discussion of ablation results and the



**Fig. 5 | Effect of MSA depth and cross-chain contacts.** **a**, Backbone accuracy (IDDT-C $\alpha$ ) for the redundancy-reduced set of the PDB after our training data cut-off, restricting to proteins in which at most 25% of the long-range contacts are between different heteromer chains. We further consider two groups of proteins based on template coverage at 30% sequence identity: covering more than 60% of the chain ( $n = 6,743$  protein chains) and covering less than 30% of the chain ( $n = 1,596$  protein chains). MSA depth is computed by counting the

number of non-gap residues for each position in the MSA (using the  $N_{\text{eff}}$  weighting scheme; see Methods for details) and taking the median across residues. The curves are obtained through Gaussian kernel average smoothing (window size is 0.2 units in  $\log_{10}(N_{\text{eff}})$ ); the shaded area is the 95% confidence interval estimated using bootstrap of 10,000 samples. **b**, An intertwined homotrimer (PDB 6SK0) is correctly predicted without input stoichiometry and only a weak template (blue is predicted and green is experimental).

effect of MSA depth on each ablation are provided in Supplementary Methods 1.13 and Supplementary Fig. 10.

### MSA depth and cross-chain contacts

Although AlphaFold has a high accuracy across the vast majority of deposited PDB structures, we note that there are still factors that affect accuracy or limit the applicability of the model. The model uses MSAs and the accuracy decreases substantially when the median alignment depth is less than around 30 sequences (see Fig. 5a for details). We observe a threshold effect where improvements in MSA depth over around 100 sequences lead to small gains. We hypothesize that the MSA information is needed to coarsely find the correct structure within the early stages of the network, but refinement of that prediction into a high-accuracy model does not depend crucially on the MSA information. The other substantial limitation that we have observed is that AlphaFold is much weaker for proteins that have few intra-chain or homotypic contacts compared to the number of heterotypic contacts (further details are provided in a companion paper<sup>39</sup>). This typically occurs for bridging domains within larger complexes in which the shape of the protein is created almost entirely by interactions with other chains in the complex. Conversely, AlphaFold is often able to give high-accuracy predictions for homomers, even when the chains are substantially intertwined (Fig. 5b). We expect that the ideas of AlphaFold are readily applicable to predicting full hetero-complexes in a future system and that this will remove the difficulty with protein chains that have a large number of hetero-contacts.

### Related work

The prediction of protein structures has had a long and varied development, which is extensively covered in a number of reviews<sup>14,40–43</sup>. Despite the long history of applying neural networks to structure prediction<sup>14,42,43</sup>, they have only recently come to improve structure prediction<sup>10,11,44,45</sup>. These approaches effectively leverage the rapid improvement in computer vision systems<sup>46</sup> by treating the problem of protein structure prediction as converting an ‘image’ of evolutionary couplings<sup>22–24</sup> to an ‘image’ of the protein distance matrix and then integrating the distance predictions into a heuristic system that produces the final 3D coordinate prediction. A few recent studies have been developed to predict the 3D coordinates directly<sup>47–50</sup>, but the accuracy of these approaches does not

match traditional, hand-crafted structure prediction pipelines<sup>51</sup>. In parallel, the success of attention-based networks for language processing<sup>52</sup> and, more recently, computer vision<sup>31,53</sup> has inspired the exploration of attention-based methods for interpreting protein sequences<sup>54–56</sup>.

### Discussion

The methodology that we have taken in designing AlphaFold is a combination of the bioinformatics and physical approaches: we use a physical and geometric inductive bias to build components that learn from PDB data with minimal imposition of handcrafted features (for example, AlphaFold builds hydrogen bonds effectively without a hydrogen bond score function). This results in a network that learns far more efficiently from the limited data in the PDB but is able to cope with the complexity and variety of structural data.

In particular, AlphaFold is able to handle missing the physical context and produce accurate models in challenging cases such as intertwined homomers or proteins that only fold in the presence of an unknown haem group. The ability to handle underspecified structural conditions is essential to learning from PDB structures as the PDB represents the full range of conditions in which structures have been solved. In general, AlphaFold is trained to produce the protein structure most likely to appear as part of a PDB structure. For example, in cases in which a particular stoichiometry, ligand or ion is predictable from the sequence alone, AlphaFold is likely to produce a structure that respects those constraints implicitly.

AlphaFold has already demonstrated its utility to the experimental community, both for molecular replacement<sup>57</sup> and for interpreting cryogenic electron microscopy maps<sup>58</sup>. Moreover, because AlphaFold outputs protein coordinates directly, AlphaFold produces predictions in graphics processing unit (GPU) minutes to GPU hours depending on the length of the protein sequence (for example, around one GPU minute per model for 384 residues; see Methods for details). This opens up the exciting possibility of predicting structures at the proteome-scale and beyond—in a companion paper<sup>39</sup>, we demonstrate the application of AlphaFold to the entire human proteome<sup>39</sup>.

The explosion in available genomic sequencing techniques and data has revolutionized bioinformatics but the intrinsic challenge of experimental structure determination has prevented a similar expansion in our structural knowledge. By developing an accurate protein structure



prediction algorithm, coupled with existing large and well-curated structure and sequence databases assembled by the experimental community, we hope to accelerate the advancement of structural bioinformatics that can keep pace with the genomics revolution. We hope that AlphaFold—and computational approaches that apply its techniques for other biophysical problems—will become essential tools of modern biology.

## Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41586-021-03819-2>.

- Thompson, M. C., Yeates, T. O. & Rodriguez, J. A. Advances in methods for atomic resolution macromolecular structure determination. *Fluorescence* **9**, 667 (2020).
- Bai, X.-C., McMullan, G. & Scheres, S. H. W. How cryo-EM is revolutionizing structural biology. *Trends Biochem. Sci.* **40**, 49–57 (2015).
- Jaskolski, M., Dauter, Z. & Wlodawer, A. A brief history of macromolecular crystallography, illustrated by a family tree and its Nobel fruits. *FEBS J.* **281**, 3985–4009 (2014).
- Wüthrich, K. The way to NMR structures of proteins. *Nat. Struct. Biol.* **8**, 923–925 (2001).
- wwPDB Consortium. Protein Data Bank: the single global archive for 3D macromolecular structure data. *Nucleic Acids Res.* **47**, D520–D528 (2018).
- Mitchell, A. L. et al. MGnify: the microbiome analysis resource in 2020. *Nucleic Acids Res.* **48**, D570–D578 (2020).
- Steinberger, M., Mirdita, M. & Söding, J. Protein-level assembly increases protein sequence recovery from metagenomic samples manifold. *Nat. Methods* **16**, 603–606 (2019).
- Dill, K. A., Ozkan, S. B., Shell, M. S. & Weikel, T. R. The protein folding problem. *Annu. Rev. Biophys.* **37**, 289–316 (2008).
- Anfinsen, C. B. Principles that govern the folding of protein chains. *Science* **181**, 223–230 (1973).
- Senior, A. W. et al. Improved protein structure prediction using potentials from deep learning. *Nature* **577**, 706–710 (2020).
- Wang, S., Sun, S., Li, Z., Zhang, R. & Xu, J. Accurate de novo prediction of protein contact map by ultra-deep learning model. *PLOS Comput. Biol.* **13**, e1005324 (2017).
- Zheng, W. et al. Deep-learning contact-map guided protein structure prediction in CASP13. *Proteins* **87**, 1149–1164 (2019).
- Abriata, L. A., Tamò, G. E. & Dal Peraro, M. A further leap of improvement in tertiary structure prediction in CASP13 prompts new routes for future assessments. *Proteins* **87**, 1100–1112 (2019).
- Pearce, R. & Zhang, Y. Deep learning techniques have significantly impacted protein structure prediction and protein design. *Curr. Opin. Struct. Biol.* **68**, 194–207 (2021).
- Moult, J., Fidelis, K., Kryshtafovych, A., Schwede, T. & Topf, M. Critical assessment of techniques for protein structure prediction, fourteenth round. *CASP 14 Abstract Book* [https://www.predictioncenter.org/casp14/doc/CASP14\\_Abstracts.pdf](https://www.predictioncenter.org/casp14/doc/CASP14_Abstracts.pdf) (2020).
- Brini, E., Simmerling, C. & Dill, K. Protein storytelling through physics. *Science* **370**, eaaz3041 (2020).
- Sippl, M. J. Calculation of conformational ensembles from potentials of mean force. An approach to the knowledge-based prediction of local structures in globular proteins. *J. Mol. Biol.* **213**, 859–883 (1990).
- Šali, A. & Blundell, T. L. Comparative protein modelling by satisfaction of spatial restraints. *J. Mol. Biol.* **234**, 779–815 (1993).
- Roy, A., Kucukural, A. & Zhang, Y. I-TASSER: a unified platform for automated protein structure and function prediction. *Nat. Protocols* **5**, 725–738 (2010).
- Altschuh, D., Lesk, A. M., Bloomer, A. C. & Klug, A. Correlation of co-ordinated amino acid substitutions with function in viruses related to tobacco mosaic virus. *J. Mol. Biol.* **193**, 693–707 (1987).
- Shindyalov, I. N., Kolchanov, N. A. & Sander, C. Can three-dimensional contacts in protein structures be predicted by analysis of correlated mutations? *Protein Eng.* **7**, 349–358 (1994).
- Weigt, M., White, R. A., Szurmant, H., Hoch, J. A. & Hwa, T. Identification of direct residue contacts in protein–protein interaction by message passing. *Proc. Natl Acad. Sci. USA* **106**, 67–72 (2009).
- Marks, D. S. et al. Protein 3D structure computed from evolutionary sequence variation. *PLoS ONE* **6**, e28766 (2011).
- Jones, D. T., Buchan, D. W. A., Cozzetto, D. & Pontil, M. PSICOV: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments. *Bioinformatics* **28**, 184–190 (2012).
- Moult, J., Pedersen, J. T., Judson, R. & Fidelis, K. A large-scale experiment to assess protein structure prediction methods. *Proteins* **23**, ii–iv (1995).
- Kryshtafovych, A., Schwede, T., Topf, M., Fidelis, K. & Moult, J. Critical assessment of methods of protein structure prediction (CASP)-round XIII. *Proteins* **87**, 1011–1020 (2019).
- Zhang, Y. & Skolnick, J. Scoring function for automated assessment of protein structure template quality. *Proteins* **57**, 702–710 (2004).
- Tu, Z. & Bai, X. Auto-context and its application to high-level vision tasks and 3D brain image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **32**, 1744–1757 (2010).
- Carreira, J., Agrawal, P., Fragkiadaki, K. & Malik, J. Human pose estimation with iterative error feedback. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition* 4733–4742 (2016).
- Mirabello, C. & Wallner, B. rawMSA: end-to-end deep learning using raw multiple sequence alignments. *PLoS ONE* **14**, e0220182 (2019).
- Huang, Z. et al. CCNet: criss-cross attention for semantic segmentation. In *Proc. IEEE/CVF International Conference on Computer Vision* 603–612 (2019).
- Hornak, V. et al. Comparison of multiple Amber force fields and development of improved protein backbone parameters. *Proteins* **65**, 712–725 (2006).
- Zemla, A. LGA: a method for finding 3D similarities in protein structures. *Nucleic Acids Res.* **31**, 3370–3374 (2003).
- Mariani, V., Biasini, M., Barbato, A. & Schwede, T. LDDT: a local superposition-free score for comparing protein structures and models using distance difference tests. *Bioinformatics* **29**, 2722–2728 (2013).
- Xie, Q., Luong, M.-T., Hovy, E. & Le, Q. V. Self-training with noisy student improves imagenet classification. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition* 10687–10698 (2020).
- Mirdita, M. et al. Uniclust databases of clustered and deeply annotated protein sequences and alignments. *Nucleic Acids Res.* **45**, D170–D176 (2017).
- Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proc. 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* 1, 4171–4186 (2019).
- Rao, R. et al. MSA transformer. In *Proc. 38th International Conference on Machine Learning* PMLR 139, 8844–8856 (2021).
- Tunyasuvunakool, K. et al. Highly accurate protein structure prediction for the human proteome. *Nature* <https://doi.org/10.1038/s41586-021-03828-1> (2021).
- Kuhlman, B. & Bradley, P. Advances in protein structure prediction and design. *Nat. Rev. Mol. Cell Biol.* **20**, 681–697 (2019).
- Marks, D. S., Hopf, T. A. & Sander, C. Protein structure prediction from sequence variation. *Nat. Biotechnol.* **30**, 1072–1080 (2012).
- Qian, N. & Sejnowski, T. J. Predicting the secondary structure of globular proteins using neural network models. *J. Mol. Biol.* **202**, 865–884 (1988).
- Fariselli, P., Olmea, O., Valencia, A. & Casadio, R. Prediction of contact maps with neural networks and correlated mutations. *Protein Eng.* **14**, 835–843 (2001).
- Yang, J. et al. Improved protein structure prediction using predicted interresidue orientations. *Proc. Natl Acad. Sci. USA* **117**, 1496–1503 (2020).
- Li, Y. et al. Deducing high-accuracy protein contact-maps from a triplet of coevolutionary matrices through deep residual convolutional networks. *PLOS Comput. Biol.* **17**, e1008865 (2021).
- He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition* 770–778 (2016).
- AlQuraishi, M. End-to-end differentiable learning of protein structure. *Cell Syst.* **8**, 292–301 (2019).
- Senior, A. W. et al. Protein structure prediction using multiple deep neural networks in the 13th Critical Assessment of Protein Structure Prediction (CASP13). *Proteins* **87**, 1141–1148 (2019).
- Ingraham, J., Riesselman, A. J., Sander, C. & Marks, D. S. Learning protein structure with a differentiable simulator. In *Proc. International Conference on Learning Representations* (2019).
- Li, J. Universal transforming geometric network. Preprint at <https://arxiv.org/abs/1908.00723> (2019).
- Xu, J., McPartlon, M. & Li, J. Improved protein structure prediction by deep learning irrespective of co-evolution information. *Nat. Mach. Intell.* **3**, 601–609 (2021).
- Vaswani, A. et al. Attention is all you need. In *Advances in Neural Information Processing Systems* 5998–6008 (2017).
- Wang, H. et al. Axial-deeplab: stand-alone axial-attention for panoptic segmentation. In *European Conference on Computer Vision* 108–126 (Springer, 2020).
- Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G. M. Unified rational protein engineering with sequence-based deep representation learning. *Nat. Methods* **16**, 1315–1322 (2019).
- Heinzinger, M. et al. Modeling aspects of the language of life through transfer-learning protein sequences. *BMC Bioinformatics* **20**, 723 (2019).
- Rives, A. et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proc. Natl Acad. Sci. USA* **118**, e2016239118 (2021).
- Pereira, J. et al. High-accuracy protein structure prediction in CASP14. *Proteins* <https://doi.org/10.1002/prot.26171> (2021).
- Gupta, M. et al. CryoEM and AI reveal a structure of SARS-CoV-2 Nsp2, a multifunctional protein involved in key host processes. Preprint at <https://doi.org/10.1101/2021.05.10.443524> (2021).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021

## Methods

### Full algorithm details

Extensive explanations of the components and their motivations are available in Supplementary Methods 1.1–1.10, in addition, pseudocode is available in Supplementary Information Algorithms 1–32, network diagrams in Supplementary Figs. 1–8, input features in Supplementary Table 1 and additional details are provided in Supplementary Tables 2, 3. Training and inference details are provided in Supplementary Methods 1.11–1.12 and Supplementary Tables 4, 5.

### IPA

The IPA module combines the pair representation, the single representation and the geometric representation to update the single representation (Supplementary Fig. 8). Each of these representations contributes affinities to the shared attention weights and then uses these weights to map its values to the output. The IPA operates in 3D space. Each residue produces query points, key points and value points in its local frame. These points are projected into the global frame using the backbone frame of the residue in which they interact with each other. The resulting points are then projected back into the local frame. The affinity computation in the 3D space uses squared distances and the coordinate transformations ensure the invariance of this module with respect to the global frame (see Supplementary Methods 1.8.2 ‘Invariant point attention (IPA)’ for the algorithm, proof of invariance and a description of the full multi-head version). A related construction that uses classic geometric invariants to construct pairwise features in place of the learned 3D points has been applied to protein design<sup>59</sup>.

In addition to the IPA, standard dot product attention is computed on the abstract single representation and a special attention on the pair representation. The pair representation augments both the logits and the values of the attention process, which is the primary way in which the pair representation controls the structure generation.

### Inputs and data sources

Inputs to the network are the primary sequence, sequences from evolutionarily related proteins in the form of a MSA created by standard tools including jackhmmer<sup>60</sup> and HHBlits<sup>61</sup>, and 3D atom coordinates of a small number of homologous structures (templates) where available. For both the MSA and templates, the search processes are tuned for high recall; spurious matches will probably appear in the raw MSA but this matches the training condition of the network.

One of the sequence databases used, Big Fantastic Database (BFD), was custom-made and released publicly (see ‘Data availability’) and was used by several CASP teams. BFD is one of the largest publicly available collections of protein families. It consists of 65,983,866 families represented as MSAs and hidden Markov models (HMMs) covering 2,204,359,010 protein sequences from reference databases, metagenomes and metatranscriptomes.

BFD was built in three steps. First, 2,423,213,294 protein sequences were collected from UniProt (Swiss-Prot&TrEMBL, 2017-11)<sup>62</sup>, a soil reference protein catalogue and the marine eukaryotic reference catalogue<sup>7</sup>, and clustered to 30% sequence identity, while enforcing a 90% alignment coverage of the shorter sequences using MMseqs2/LinClust<sup>63</sup>. This resulted in 345,159,030 clusters. For computational efficiency, we removed all clusters with less than three members, resulting in 61,083,719 clusters. Second, we added 166,510,624 representative protein sequences from Metaclust NR (2017-05; discarding all sequences shorter than 150 residues)<sup>63</sup> by aligning them against the cluster representatives using MMseqs2<sup>64</sup>. Sequences that fulfilled the sequence identity and coverage criteria were assigned to the best scoring cluster. The remaining 25,347,429 sequences that could not be assigned were clustered separately and added as new clusters, resulting in the final clustering. Third, for each of the clusters, we computed an MSA using

FAMSA<sup>65</sup> and computed the HMMs following the Uniclust HH-suite database protocol<sup>36</sup>.

The following versions of public datasets were used in this study. Our models were trained on a copy of the PDB<sup>5</sup> downloaded on 28 August 2019. For finding template structures at prediction time, we used a copy of the PDB downloaded on 14 May 2020, and the PDB70<sup>66</sup> clustering database downloaded on 13 May 2020. For MSA lookup at both training and prediction time, we used Uniref90<sup>67</sup> v.2020\_01, BFD, Uniclust30<sup>36</sup> v.2018\_08 and Mgnify<sup>6</sup> v.2018\_12. For sequence distillation, we used Uniclust30<sup>36</sup> v.2018\_08 to construct a distillation structure dataset. Full details are provided in Supplementary Methods 1.2.

For MSA search on BFD + Uniclust30, and template search against PDB70, we used HHBlits<sup>61</sup> and HHSearch<sup>66</sup> from hh-suite v.3.0-beta.3 (version 14/07/2017). For MSA search on Uniref90 and clustered Mgnify, we used jackhmmer from HMMER3<sup>68</sup>. For constrained relaxation of structures, we used OpenMM v.7.3.1<sup>69</sup> with the Amber99sb force field<sup>32</sup>. For neural network construction, running and other analyses, we used TensorFlow<sup>70</sup>, Sonnet<sup>71</sup>, NumPy<sup>72</sup>, Python<sup>73</sup> and Colab<sup>74</sup>.

To quantify the effect of the different sequence data sources, we re-ran the CASP14 proteins using the same models but varying how the MSA was constructed. Removing BFD reduced the mean accuracy by 0.4 GDT, removing Mgnify reduced the mean accuracy by 0.7 GDT, and removing both reduced the mean accuracy by 6.1 GDT. In each case, we found that most targets had very small changes in accuracy but a few outliers had very large (20+ GDT) differences. This is consistent with the results in Fig. 5a in which the depth of the MSA is relatively unimportant until it approaches a threshold value of around 30 sequences when the MSA size effects become quite large. We observe mostly overlapping effects between inclusion of BFD and Mgnify, but having at least one of these metagenomics databases is very important for target classes that are poorly represented in UniRef, and having both was necessary to achieve full CASP accuracy.

### Training regimen

To train, we use structures from the PDB with a maximum release date of 30 April 2018. Chains are sampled in inverse proportion to cluster size of a 40% sequence identity clustering. We then randomly crop them to 256 residues and assemble into batches of size 128. We train the model on Tensor Processing Unit (TPU) v3 with a batch size of 1 per TPU core, hence the model uses 128 TPU v3 cores. The model is trained until convergence (around 10 million samples) and further fine-tuned using longer crops of 384 residues, larger MSA stack and reduced learning rate (see Supplementary Methods 1.11 for the exact configuration). The initial training stage takes approximately 1 week, and the fine-tuning stage takes approximately 4 additional days.

The network is supervised by the FAPE loss and a number of auxiliary losses. First, the final pair representation is linearly projected to a binned distance distribution (distogram) prediction, scored with a cross-entropy loss. Second, we use random masking on the input MSAs and require the network to reconstruct the masked regions from the output MSA representation using a BERT-like loss<sup>37</sup>. Third, the output single representations of the structure module are used to predict binned per-residue IDDT-C $\alpha$  values. Finally, we use an auxiliary side-chain loss during training, and an auxiliary structure violation loss during fine-tuning. Detailed descriptions and weighting are provided in the Supplementary Information.

An initial model trained with the above objectives was used to make structure predictions for a Uniclust dataset of 355,993 sequences with the full MSAs. These predictions were then used to train a final model with identical hyperparameters, except for sampling examples 75% of the time from the Uniclust prediction set, with sub-sampled MSAs, and 25% of the time from the clustered PDB set.

We train five different models using different random seeds, some with templates and some without, to encourage diversity in the predictions (see Supplementary Table 5 and Supplementary Methods 1.12.1



for details). We also fine-tuned these models after CASP14 to add a pTM prediction objective (Supplementary Methods 1.9.7) and use the obtained models for Fig. 2d.

### Inference regimen

We inference the five trained models and use the predicted confidence score to select the best model per target.

Using our CASP14 configuration for AlphaFold, the trunk of the network is run multiple times with different random choices for the MSA cluster centres (see Supplementary Methods 1.11.2 for details of the ensembling procedure). The full time to make a structure prediction varies considerably depending on the length of the protein. Representative timings for the neural network using a single model on V100 GPU are 4.8 min with 256 residues, 9.2 min with 384 residues and 18 h at 2,500 residues. These timings are measured using our open-source code, and the open-source code is notably faster than the version we ran in CASP14 as we now use the XLA compiler<sup>75</sup>.

Since CASP14, we have found that the accuracy of the network without ensembling is very close or equal to the accuracy with ensembling and we turn off ensembling for most inference. Without ensembling, the network is 8× faster and the representative timings for a single model are 0.6 min with 256 residues, 1.1 min with 384 residues and 2.1 h with 2,500 residues.

Inferencing large proteins can easily exceed the memory of a single GPU. For a V100 with 16 GB of memory, we can predict the structure of proteins up to around 1,300 residues without ensembling and the 256- and 384-residue inference times are using the memory of a single GPU. The memory usage is approximately quadratic in the number of residues, so a 2,500-residue protein involves using unified memory so that we can greatly exceed the memory of a single V100. In our cloud setup, a single V100 is used for computation on a 2,500-residue protein but we requested four GPUs to have sufficient memory.

Searching genetic sequence databases to prepare inputs and final relaxation of the structures take additional central processing unit (CPU) time but do not require a GPU or TPU.

### Metrics

The predicted structure is compared to the true structure from the PDB in terms of IDDT metric<sup>34</sup>, as this metric reports the domain accuracy without requiring a domain segmentation of chain structures. The distances are either computed between all heavy atoms (IDDT) or only the C $\alpha$  atoms to measure the backbone accuracy (IDDT-C $\alpha$ ). As IDDT-C $\alpha$  only focuses on the C $\alpha$  atoms, it does not include the penalty for structural violations and clashes. Domain accuracies in CASP are reported as GDT<sup>33</sup> and the TM-score<sup>27</sup> is used as a full chain global superposition metric.

We also report accuracies using the r.m.s.d.<sub>95</sub> (C $\alpha$  r.m.s.d. at 95% coverage). We perform five iterations of (1) a least-squares alignment of the predicted structure and the PDB structure on the currently chosen C $\alpha$  atoms (using all C $\alpha$  atoms in the first iteration); (2) selecting the 95% of C $\alpha$  atoms with the lowest alignment error. The r.m.s.d. of the atoms chosen for the final iterations is the r.m.s.d.<sub>95</sub>. This metric is more robust to apparent errors that can originate from crystal structure artefacts, although in some cases the removed 5% of residues will contain genuine modelling errors.

### Test set of recent PDB sequences

For evaluation on recent PDB sequences (Figs. 2a–d, 4a, 5a), we used a copy of the PDB downloaded 15 February 2021. Structures were filtered to those with a release date after 30 April 2018 (the date limit for inclusion in the training set for AlphaFold). Chains were further filtered to remove sequences that consisted of a single amino acid as well as sequences with an ambiguous chemical component at any residue position. Exact duplicates were removed, with the chain with the most resolved C $\alpha$  atoms used as the representative sequence. Subsequently,

structures with less than 16 resolved residues, with unknown residues or solved by NMR methods were removed. As the PDB contains many near-duplicate sequences, the chain with the highest resolution was selected from each cluster in the PDB 40% sequence clustering of the data. Furthermore, we removed all sequences for which fewer than 80 amino acids had the alpha carbon resolved and removed chains with more than 1,400 residues. The final dataset contained 10,795 protein sequences.

The procedure for filtering the recent PDB dataset based on prior template identity was as follows. Hmsearch was run with default parameters against a copy of the PDB SEQRES fasta downloaded 15 February 2021. Template hits were accepted if the associated structure had a release date earlier than 30 April 2018. Each residue position in a query sequence was assigned the maximum identity of any template hit covering that position. Filtering then proceeded as described in the individual figure legends, based on a combination of maximum identity and sequence coverage.

The MSA depth analysis was based on computing the normalized number of effective sequences ( $N_{\text{eff}}$ ) for each position of a query sequence. Per-residue  $N_{\text{eff}}$  values were obtained by counting the number of non-gap residues in the MSA for this position and weighting the sequences using the  $N_{\text{eff}}$  scheme<sup>76</sup> with a threshold of 80% sequence identity measured on the region that is non-gap in either sequence.

### Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this paper.

### Data availability

All input data are freely available from public sources.

Structures from the PDB were used for training and as templates (<https://www.wwpdb.org/ftp/pdb-ftp-sites>; for the associated sequence data and 40% sequence clustering see also [https://ftp.wwpdb.org/pub/pdb/derived\\_data/](https://ftp.wwpdb.org/pub/pdb/derived_data/) and <https://cdn.rcsb.org/resources/sequence/clusters/bc-40.out>). Training used a version of the PDB downloaded 28 August 2019, while the CASP14 template search used a version downloaded 14 May 2020. The template search also used the PDB70 database, downloaded 13 May 2020 ([https://wwwuser.gwdg.de/~compbiol/data/hhsuite/databases/hhsuite\\_dbs/](https://wwwuser.gwdg.de/~compbiol/data/hhsuite/databases/hhsuite_dbs/)).

We show experimental structures from the PDB with accession numbers 6Y4F<sup>77</sup>, 6YJ1<sup>78</sup>, 6VR4<sup>79</sup>, 6SK0<sup>80</sup>, 6FES<sup>81</sup>, 6W6W<sup>82</sup>, 6T1Z<sup>83</sup> and 7JTL<sup>84</sup>.

For MSA lookup at both the training and prediction time, we used UniRef90 v.2020\_01 ([https://ftp.ebi.ac.uk/pub/databases/uniprot/previous\\_releases/release-2020\\_01/uniref/](https://ftp.ebi.ac.uk/pub/databases/uniprot/previous_releases/release-2020_01/uniref/)), BFD (<https://bfd.mmseqs.com>), Uniclust30 v.2018\_08 ([https://wwwuser.gwdg.de/~compbiol/uniclust/2018\\_08/](https://wwwuser.gwdg.de/~compbiol/uniclust/2018_08/)) and MGnify clusters v.2018\_12 ([https://ftp.ebi.ac.uk/pub/databases/metagenomics/peptide\\_database/2018\\_12/](https://ftp.ebi.ac.uk/pub/databases/metagenomics/peptide_database/2018_12/)). Uniclust30 v.2018\_08 was also used as input for constructing a distillation structure dataset.

### Code availability

Source code for the AlphaFold model, trained weights and inference script are available under an open-source license at <https://github.com/deepmind/alphafold>.

Neural networks were developed with TensorFlow v.1 (<https://github.com/tensorflow/tensorflow>), Sonnet v.1 (<https://github.com/deepmind/sonnet>), JAX v.0.1.69 (<https://github.com/google/jax/>) and Haiku v.0.0.4 (<https://github.com/deepmind/dm-haiku>). The XLA compiler is bundled with JAX and does not have a separate version number.

For MSA search on BFD+Uniclust30, and for template search against PDB70, we used HHBlits and HHSearch from hh-suite v.3.0-beta.3 release 14/07/2017 (<https://github.com/soedinglab/hh-suite>). For MSA search on UniRef90 and clustered MGnify, we used jackhmmer from

HMMER v.3.3 (<http://eddylab.org/software/hmmer/>). For constrained relaxation of structures, we used OpenMM v.7.3.1 (<https://github.com/openmm/openmm>) with the Amber99sb force field.

Construction of BFD used MMseqs2 v.925AF (<https://github.com/soedinglab/MMseqs2>) and FAMSA v.1.2.5 (<https://github.com/refresh-bio/FAMSA>).

Data analysis used Python v.3.6 (<https://www.python.org/>), NumPy v.1.16.4 (<https://github.com/numpy/numpy>), SciPy v.1.2.1 (<https://www.scipy.org/>), seaborn v.0.11.1 (<https://github.com/mwaskom/seaborn>), Matplotlib v.3.3.4 (<https://github.com/matplotlib/matplotlib>), bokeh v.1.4.0 (<https://github.com/bokeh/bokeh>), pandas v.1.1.5 (<https://github.com/pandas-dev/pandas>), plotnine v.0.8.0 (<https://github.com/has2k1/plotnine>), statsmodels v.0.12.2 (<https://github.com/statsmodels/statsmodels>) and Colab (<https://research.google.com/colaboratory>). TM-align v.20190822 (<https://zhanglab.dcm.med.umich.edu/TM-align/>) was used for computing TM-scores. Structure visualizations were created in Pymol v.2.3.0 (<https://github.com/schrodinger/pymol-open-source>).

59. Ingraham, J., Garg, V. K., Barzilay, R. & Jaakkola, T. Generative models for graph-based protein design. in *Proc. 33rd Conference on Neural Information Processing Systems* (2019).
60. Johnson, L. S., Eddy, S. R. & Portugaly, E. Hidden Markov model speed heuristic and iterative HMM search procedure. *BMC Bioinformatics* **11**, 431 (2010).
61. Remmert, M., Biegert, A., Hauser, A. & Söding, J. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nat. Methods* **9**, 173–175 (2012).
62. The UniProt Consortium. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Res.* **49**, D480–D489 (2020).
63. Steinegger, M. & Söding, J. Clustering huge protein sequence sets in linear time. *Nat. Commun.* **9**, 2542 (2018).
64. Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat. Biotechnol.* **35**, 1026–1028 (2017).
65. Deorowicz, S., Debudaj-Grabysz, A. & Gudyś, A. FAMSA: fast and accurate multiple sequence alignment of huge protein families. *Sci. Rep.* **6**, 33964 (2016).
66. Steinegger, M. et al. HH-suite3 for fast remote homology detection and deep protein annotation. *BMC Bioinformatics* **20**, 473 (2019).
67. Suze, B. E., Wang, Y., Huang, H., McGarvey, P. B. & Wu, C. H. UniRef clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics* **31**, 926–932 (2015).
68. Eddy, S. R. Accelerated profile HMM searches. *PLOS Comput. Biol.* **7**, e1002195 (2011).
69. Eastman, P. et al. OpenMM 7: rapid development of high performance algorithms for molecular dynamics. *PLOS Comput. Biol.* **13**, e1005659 (2017).
70. Ashish, A. M. A. et al. TensorFlow: large-scale machine learning on heterogeneous systems. Preprint at <https://arxiv.org/abs/1603.04467> (2015).
71. Reynolds, M. et al. Open sourcing Sonnet – a new library for constructing neural networks. *DeepMind* <https://deepmind.com/blog/open-sourcing-sonnet/> (7 April 2017).
72. Harris, C. R. et al. Array programming with NumPy. *Nature* **585**, 357–362 (2020).
73. Van Rossum, G. & Drake, F. L. *Python 3 Reference Manual* (CreateSpace, 2009).
74. Bisong, E. in *Building Machine Learning and Deep Learning Models on Google Cloud Platform: A Comprehensive Guide for Beginners* 59–64 (Apress, 2019).
75. TensorFlow. XLA: Optimizing Compiler for TensorFlow. <https://www.tensorflow.org/xla> (2018).

76. Wu, T., Hou, J., Adhikari, B. & Cheng, J. Analysis of several key factors influencing deep learning-based inter-residue contact prediction. *Bioinformatics* **36**, 1091–1098 (2020).
77. Jiang, W. et al. MrpH, a new class of metal-binding adhesin, requires zinc to mediate biofilm formation. *PLoS Pathog.* **16**, e1008707 (2020).
78. Dunne, M., Ernst, P., Sobieraj, A., Pluckthun, A. & Loessner, M. J. The M23 peptidase domain of the Staphylococcal phage 2638A endolysin. *PDB* <https://doi.org/10.2210/pdb6YJ1/pdb> (2020).
79. Drobysheva, A. V. et al. Structure and function of virion RNA polymerase of a crAss-like phage. *Nature* **589**, 306–309 (2021).
80. Flaughnatti, N. et al. Structural basis for loading and inhibition of a bacterial T6SS phospholipase effector by the VgrG spike. *EMBO J.* **39**, e104129 (2020).
81. ElGamacy, M. et al. An interface-driven design strategy yields a novel, corrugated protein architecture. *ACS Synth. Biol.* **7**, 2226–2235 (2018).
82. Lim, C. J. et al. The structure of human CST reveals a decameric assembly bound to telomeric DNA. *Science* **368**, 1081–1085 (2020).
83. Debruycker, V. et al. An embedded lipid in the multidrug transporter LmrP suggests a mechanism for polyspecificity. *Nat. Struct. Mol. Biol.* **27**, 829–835 (2020).
84. Flower, T. G. et al. Structure of SARS-CoV-2 ORF8, a rapidly evolving immune evasion protein. *Proc. Natl Acad. Sci. USA* **118**, e2021785118 (2021).

**Acknowledgements** We thank A. Rustemi, A. Gu, A. Guseynov, B. Hechtman, C. Beattie, C. Jones, C. Donner, E. Parisotto, E. Elsen, F. Popovici, G. Necula, H. Maclean, J. Menick, J. Kirkpatrick, J. Molloy, J. Yim, J. Stanway, K. Simonyan, L. Sifre, L. Martens, M. Johnson, M. O’Neill, N. Antropova, R. Hadsell, S. Blackwell, S. Das, S. Hou, S. Gouws, S. Wheelwright, T. Hennigan, T. Ward, Z. Wu, Z. Wu, Z. Avsec and the Research Platform Team for their contributions; M. Mirdita for his help with the datasets; M. Piovesan-Forster, A. Nelson and R. Kemp for their help managing the project; the JAX, TensorFlow and XLA teams for detailed support and enabling machine learning models of the complexity of AlphaFold; our colleagues at DeepMind, Google and Alphabet for their encouragement and support; and J. Moult and the CASP14 organizers, and the experimentalists whose structures enabled the assessment. M.S. acknowledges support from the National Research Foundation of Korea grant (2019R1A6A1A10073437, 2020M3A9G7I03933) and the Creative-Pioneering Researchers Program through Seoul National University.

**Author contributions** J.J. and D.H. led the research. J.J., R.E., A. Pritzel, M.F., O.R., R.B., A. Potapenko, S.A.A.K., B.R.-P., J.A., M.P., T. Berghammer and O.V. developed the neural network architecture and training. T.G., A.Ž., K.T., R.B., A.B., R.E., A.J.B., A.C., S.N., R.J., D.R., M.Z. and S.B. developed the data, analytics and inference systems. D.H., K.K., P.K., C.M. and E.C. managed the research. T.G. led the technical platform. P.K., A.W.S., K.K., O.V., D.S., S.P. and T. Back contributed technical advice and ideas. M.S. created the BFD genomics database and provided technical assistance on HHblits. D.H., R.E., A.W.S. and K.K. conceived the AlphaFold project. J.J., R.E. and A.W.S. conceived the end-to-end approach. J.J., A. Pritzel, O.R., A. Potapenko, R.E., M.F., T.G., K.T., C.M. and D.H. wrote the paper.

**Competing interests** J.J., R.E., A. Pritzel, T.G., M.F., O.R., R.B., A.B., S.A.A.K., D.R. and A.W.S. have filed non-provisional patent applications 16/701,070 and PCT/EP2020/084238, and provisional patent applications 63/107,362, 63/118,917, 63/118,918, 63/118,921 and 63/118,919, each in the name of DeepMind Technologies Limited, each pending, relating to machine learning for predicting protein structures. The other authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41586-021-03819-2>.

**Correspondence and requests for materials** should be addressed to J.J. or D.H.

**Peer review information** Nature thanks Mohammed AlQuraishi, Charlotte Deane and Yang Zhang for their contribution to the peer review of this work.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>.

## Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

### Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- ☐ ☒ The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement
- ☒ ☐ A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- ☐ ☒ The statistical test(s) used AND whether they are one- or two-sided  
*Only common tests should be described solely by name; describe more complex techniques in the Methods section.*
- ☒ ☐ A description of all covariates tested
- ☐ ☒ A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- ☐ ☒ A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- ☒ ☐ For null hypothesis testing, the test statistic (e.g.  $F$ ,  $t$ ,  $r$ ) with confidence intervals, effect sizes, degrees of freedom and  $P$  value noted  
*Give  $P$  values as exact values whenever suitable.*
- ☐ ☒ For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- ☒ ☐ For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- ☐ ☒ Estimates of effect sizes (e.g. Cohen's  $d$ , Pearson's  $r$ ), indicating how they were calculated

*Our web collection on [statistics for biologists](#) contains articles on many of the points above.*

### Software and code

Policy information about [availability of computer code](#)

#### Data collection

Source code for the AlphaFold model, trained weights, and inference script will be made available under an open-source license at <https://github.com/deepmind/> upon publication.

Neural networks were developed with TensorFlow v1 (<https://github.com/tensorflow/tensorflow>), Sonnet v1 (<https://github.com/deepmind/sonnet>), JAX v0.1.69 (<https://github.com/google/jax>), and Haiku v0.0.4 (<https://github.com/deepmind/dm-haiku>). The XLA compiler is bundled with JAX and does not have a separate version number.

For MSA search on BFD+Uniclust30, and for template search against PDB70, we used HHBlits and HHSearch from hh-suite v3.0-beta.3 14/07/2017 (<https://github.com/soedinglab/hh-suite>). For MSA search on UniRef90 and clustered MGnify, we used jackhmmer from HMMER v3.3 (<http://eddylib.org/software/hmmer/>). For constrained relaxation of structures, we used OpenMM v7.3.1 (<https://github.com/openmm/openmm>) with the Amber99sb force field.

Construction of BFD used MMseqs2 version 925AF (<https://github.com/soedinglab/MMseqs2>) and FAMSA v1.2.5 (<https://github.com/refresh-bio/FAMSA>).

#### Data analysis

Data analysis used Python v3.6 (<https://www.python.org/>), NumPy v1.16.4 (<https://github.com/numpy/numpy>), SciPy v1.2.1 (<https://www.scipy.org/>), seaborn v0.11.1 (<https://github.com/mwaskom/seaborn>), Matplotlib v3.3.4 (<https://github.com/matplotlib/matplotlib>), bokeh v1.4.0 (<https://github.com/bokeh/bokeh>), pandas v1.1.5 (<https://github.com/pandas-dev/pandas>), plotnine v0.8.0 (<https://github.com/has2k1/plotnine>), statsmodels v0.12.2 (<https://github.com/statsmodels/statsmodels>) and Colab (<https://research.google.com/colaboratory>). TM-align v20190822 (<https://zhanglab.dcm.med.umich.edu/TM-align/>) was used for computing TM-scores. Structure visualizations were created in Pymol v2.3.0 (<https://github.com/schrodinger/pymol-open-source>).

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.



## Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

All input data are freely available from public sources.

Structures from the PDB were used for training and as templates (<https://www.wwpdb.org/ftp/pdb-ftp-sites>; for the associated sequence data and 40% sequence clustering see also [https://ftp.wwpdb.org/pub/pdb/derived\\_data/](https://ftp.wwpdb.org/pub/pdb/derived_data/) and <https://cdn.rcsb.org/resources/sequence/clusters/bc-40.out>). Training used a version of the PDB downloaded 28/08/2019, while CASP14 template search used a version downloaded 14/05/2020. Template search also used the PDB70 database, downloaded 13/05/2020 ([https://wwwuser.gwdg.de/~compbiol/data/hhsuite/databases/hhsuite\\_dbs/](https://wwwuser.gwdg.de/~compbiol/data/hhsuite/databases/hhsuite_dbs/)).

We show experimental structures from the PDB with accessions 6Y4F77, 6YJ178, 6VR479, 6SK080, 6FES81, 6W6W82, 6T1Z83, and 7JTL84.

For MSA lookup at both training and prediction time, we used UniRef90 v2020\_01 ([https://ftp.ebi.ac.uk/pub/databases/uniprot/previous\\_releases/release-2020\\_01/uniref/](https://ftp.ebi.ac.uk/pub/databases/uniprot/previous_releases/release-2020_01/uniref/)), BFD (<https://bfd.mmseqs.com>), Uniclust30 v2018\_08 ([https://wwwuser.gwdg.de/~compbiol/uniclust/2018\\_08/](https://wwwuser.gwdg.de/~compbiol/uniclust/2018_08/)), and MGnify clusters v2018\_12 ([https://ftp.ebi.ac.uk/pub/databases/metagenomics/peptide\\_database/2018\\_12/](https://ftp.ebi.ac.uk/pub/databases/metagenomics/peptide_database/2018_12/)). Uniclust30 v2018\_08 was further used as input for constructing a distillation structure dataset.

## Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences ☐ Behavioural & social sciences ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

## Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	No sample size was chosen; the method was evaluated on the full CASP14 benchmark set, and all PDB chains not in the training set (subject to the exclusions noted below).
Data exclusions	The recent PDB set was filtered (see Methods for full details). Briefly this excludes chains with too few resolved residues, longer than 1400 residues, solved by NMR or with unknown/ambiguous residues. This set was also redundancy reduced (by taking representatives from a sequence clustering), and for some figures a sequence similarity-based filter was applied to remove entries too similar to the training set (see Methods and figure legends for details).
Replication	Not applicable, no experimental work is described in this study. The results are the output of a computational method which will be made available.
Randomization	Not applicable, we are not making a comparison between two groups
Blinding	Not applicable, we are not making a comparison between two groups

## Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

### Materials & experimental systems

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

### Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging