

¹ AlphaFold: Improved protein structure prediction using potentials from deep learning

³ Andrew W. Senior^{1*}, Richard Evans^{1*}, John Jumper^{1*}, James Kirkpatrick^{1*}, Laurent Sifre^{1*}, Tim Green¹,
⁴ Chongli Qin¹, Augustin Žídek¹, Alexander W. R. Nelson¹, Alex Bridgland¹, Hugo Penedones¹,
⁵ Stig Petersen¹, Karen Simonyan¹, Steve Crossan¹, Pushmeet Kohli¹, David T. Jones^{2,3}, David Silver¹,
⁶ Koray Kavukcuoglu¹, Demis Hassabis¹

⁷ ¹DeepMind, London, UK

⁸ ²The Francis Crick Institute, London, UK

⁹ ³University College London, London, UK

¹⁰ *These authors contributed equally to this work.

¹¹ Protein structure prediction aims to determine the three-dimensional shape of a protein from
¹² its amino acid sequence¹. This problem is of fundamental importance to biology as the struc-
¹³ ture of a protein largely determines its function² but can be hard to determine experimen-
¹⁴ tally. In recent years, considerable progress has been made by leveraging genetic informa-
¹⁵ tion: analysing the co-variation of homologous sequences can allow one to infer which amino
¹⁶ acid residues are in contact, which in turn can aid structure prediction³. In this work, we
¹⁷ show that we can train a neural network to accurately predict the distances between pairs
¹⁸ of residues in a protein which convey more about structure than contact predictions. With
¹⁹ this information we construct a potential of mean force⁴ that can accurately describe the
²⁰ shape of a protein. We find that the resulting potential can be optimised by a simple gradient
²¹ descent algorithm, to realise structures without the need for complex sampling procedures.
²² The resulting system, named AlphaFold, has been shown to achieve high accuracy, even for
²³ sequences with relatively few homologous sequences. In the most recent Critical Assessment
²⁴ of Protein Structure Prediction⁵ (CASP13), a blind assessment of the state of the field of pro-
²⁵ tein structure prediction, AlphaFold created high-accuracy structures (with TM-scores[†] of
²⁶ 0.7 or higher) for 24 out of 43 free modelling domains whereas the next best method, using
²⁷ sampling and contact information, achieved such accuracy for only 14 out of 43 domains.
²⁸ AlphaFold represents a significant advance in protein structure prediction. We expect the in-
²⁹ creased accuracy of structure predictions for proteins to enable insights in understanding the
³⁰ function and malfunction of these proteins, especially in cases where no homologous proteins
³¹ have been experimentally determined⁷.

³² Proteins are at the core of most biological processes. Since the function of a protein is
³³ dependent on its structure, understanding protein structure has been a grand challenge in biology
³⁴ for decades. While several experimental structure determination techniques have been developed

[†]Template Modelling score⁶, between 0 and 1, measures the degree of match of the overall (backbone) shape of a proposed structure to a native structure.

35 and improved in accuracy, they remain difficult and time-consuming². As a result, decades of
 36 theoretical work has attempted to predict protein structure from amino acid sequences.

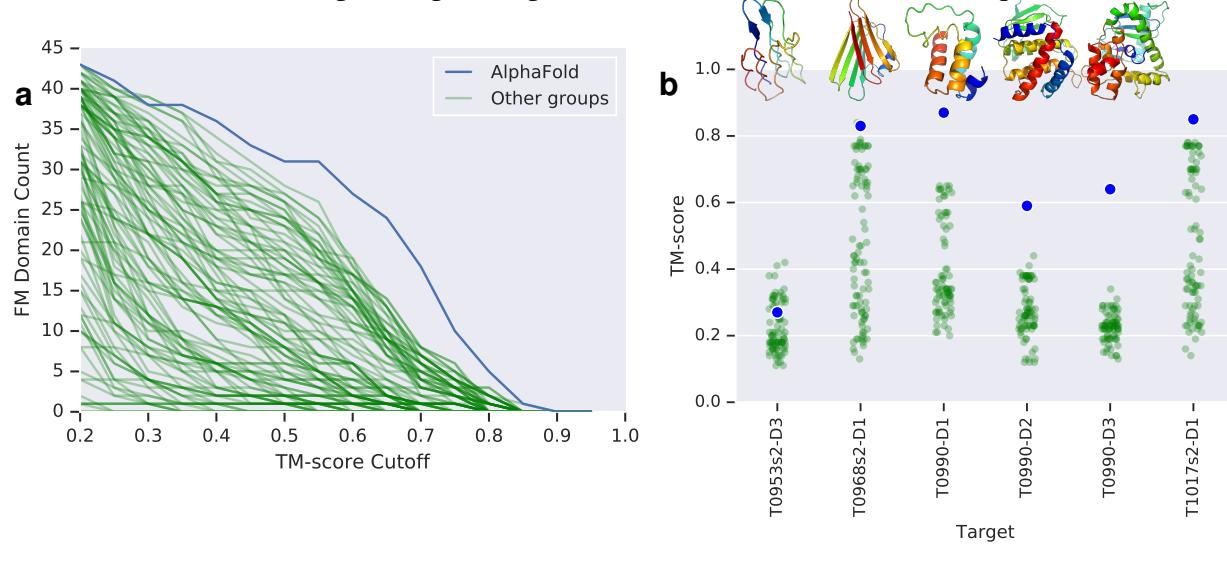


Fig. 1 | AlphaFold’s performance in the CASP13 assessment. (a) Number of free modelling (FM + FM/TBM) domains predicted to a given TM-score threshold for AlphaFold and the other 97 groups. (b) For the six new folds identified by the CASP13 assessors, AlphaFold’s TM-score compared with the other groups, with native structures. The structure of T1017s2-D1 is unavailable for publication. (c) Precisions for long-range contact prediction in CASP13 for the most probable L , $L/2$ or $L/5$ contacts, where L is the length of the domain. The distance distributions used by AlphaFold (AF) in CASP13, thresholded to contact predictions, are compared with submissions by the two best-ranked contact prediction methods in CASP13: 498 (RaptorX-Contact⁸) and 032 (TripletRes⁹), on “all groups” targets, excluding T0999.

37 CASP⁵ is a biennial blind protein structure prediction assessment run by the structure pre-
 38 diction community to benchmark progress in accuracy. In 2018, AlphaFold joined 97 groups from
 39 around the world in entering CASP13. Each group submitted up to 5 structure predictions for
 40 each of 84 protein sequences whose experimentally-determined structures were sequestered. As-
 41 sessors divided the proteins into 104 domains for scoring and classified each as being amenable
 42 to *template-based modelling* (TBM, where a protein with a similar sequence has a known struc-
 43 ture, and that homologous structure is modified in accordance with the sequence differences) or
 44 requiring *free modelling* (FM, when no homologous structure is available), with an intermediate
 45 (FM/TBM) category. Figure 1a shows that AlphaFold stands out in performance above the other
 46 entrants, predicting more FM domains to high accuracy than any other system, particularly in the

47 0.6–0.7 TM-score range. The assessors ranked the 98 participating groups by the summed, capped
48 z-scores of the structures, separated according to category. AlphaFold achieved a summed z-score
49 of 52.8 in the FM category (best-of-5) vs 36.6 for the next closest group (322)[‡]. Combining FM
50 and TBM/FM categories, AlphaFold scored 68.3 vs 48.2. AlphaFold is able to predict previously
51 unknown folds to high accuracy as shown in Figure 1b. Despite using only free modelling tech-
52 niques and not using templates, AlphaFold also scored well in the TBM category according to the
53 assessors’ formula 0-capped z-score, ranking fourth by the top-1 model or first by the best-of-5
54 models. Much of the accuracy of AlphaFold is due to the accuracy of the distance predictions,
55 which is evident from the high precision of the contact predictions of Table 1c.

56 The most successful free modelling approaches so far^{10–12} have relied on *fragment assembly*
57 to determine the shape of the protein of interest. In these approaches a structure is created through
58 a stochastic sampling process, such as simulated annealing¹³, that minimises a statistical potential
59 derived from summary statistics extracted from structures in the Protein Data Bank (PDB¹⁴). In
60 fragment assembly, a structure hypothesis is repeatedly modified, typically by changing the shape
61 of a short section, retaining changes which lower the potential, ultimately leading to low potential
62 structures. Simulated annealing requires many thousands of such moves and must be repeated
63 many times to have good coverage of low-potential structures.

64 In recent years, structure prediction accuracy has improved through the use of evolutionary
65 covariation data¹⁵ found in sets of related sequences. Sequences similar to the target sequence
66 are found by searching large datasets of protein sequences derived from DNA sequencing and
67 aligned to the target sequence to make a *multiple sequence alignment* (MSA). Correlated changes
68 in two amino acid residue positions across the sequences of the MSA can be used to infer which
69 residues might be in contact. Contacts are typically defined to occur when the β -carbon atoms of
70 two residues are within 8 Ångström of one another. Several methods have been used to predict
71 the probability that a pair of residues is in contact based on features computed from MSAs^{16–19}
72 including neural networks^{20–23}. Contact predictions are incorporated in structure prediction by
73 modifying the statistical potential to guide the folding process to structures that satisfy more of the
74 predicted contacts^{12,24}. Previous work^{25,26} has made predictions of the distance between residues,
75 particularly for distance geometry approaches^{8,27–29}. Neural network distance predictions without
76 covariation features were used to make the EPAD potential²⁶ which was used for ranking struc-
77 ture hypotheses and the QUARK pipeline¹² used a template-based distance profile restraint for
78 template-based modelling.

79 In this work we present a new, deep-learning, approach to protein structure prediction, whose
80 stages are illustrated in Figure 2a. We show that it is possible to construct a learned, protein-specific
81 potential by training a neural network (Fig. 2b) to make accurate predictions about the structure
82 of the protein given its sequence, and to predict the structure itself accurately by minimising the

[‡]Results from http://predictioncenter.org/casp13/zscores_final.cgi?formula=assessors

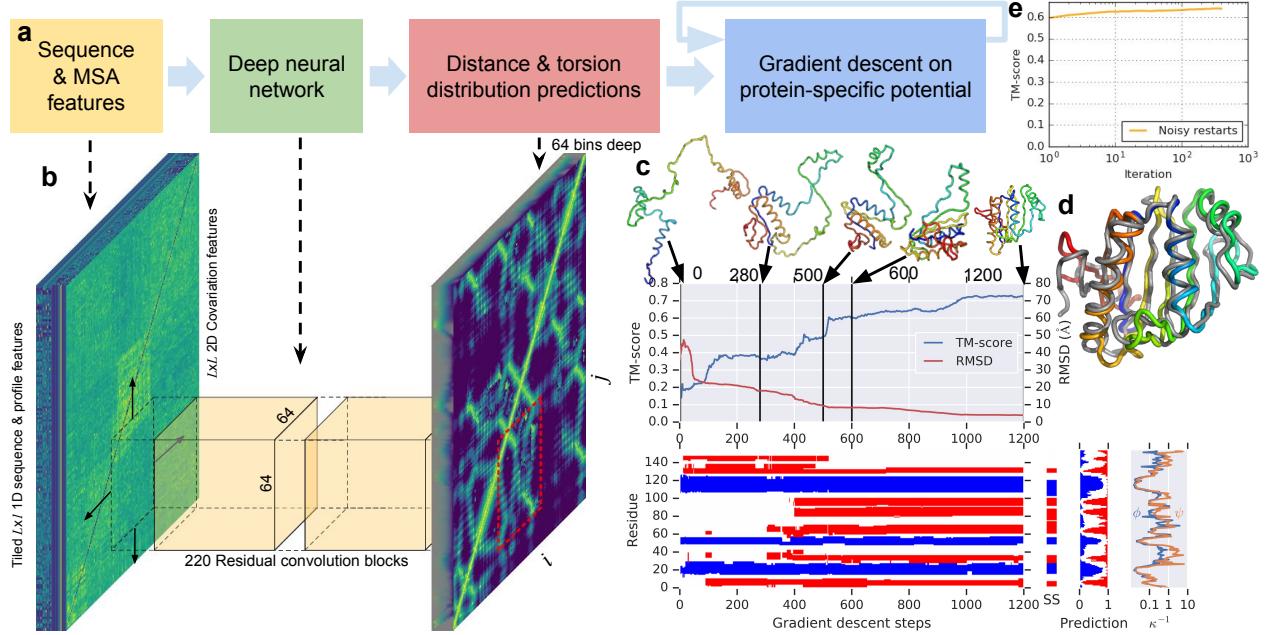


Fig. 2 | The folding process illustrated for CASP13 target T0986s2. (Length $L = 155$) (a) Steps of structure prediction. (b) The neural network predicts the entire $L \times L$ distogram based on MSA features, accumulating separate predictions for 64×64 -residue regions. (c) One iteration of gradient descent (1 200 steps) is shown, with TM-score and RMSD plotted against step number with five snapshots of the structure. The secondary structure (from SST³⁰) is also shown (helix in blue, strand in red) along with the the native secondary structure (SS), the network's secondary structure prediction probabilities and the uncertainty in torsion angle predictions (as κ^{-1} of the von Mises distributions fitted to the predictions for ϕ and ψ). While each step of gradient descent greedily lowers the potential, large global conformation changes are effected, resulting in a well-packed chain. (d) shows the final first submission overlaid on the native structure (in grey). (e) shows the average (across the test set, $n = 377$) TM-score of the lowest-potential structure against the number of repeats of gradient descent (log scale).

83 potential by gradient descent (Fig. 2c). The neural network predictions include backbone torsion
84 angles and pairwise distances between residues. Distance predictions provide more specific in-
85 formation about the structure than contact predictions and provide a richer training signal for the
86 neural network. Predicting distances, rather than contacts as in most prior work, models detailed
87 interactions rather than simple binary decisions. By jointly predicting many distances, the network
88 can propagate distance information respecting covariation, local structure and residue identities to
89 nearby residues. The predicted probability distributions can be combined to form a simple, prin-
90 cipled protein-specific potential. We show that with gradient descent, it is simple to find a set of
91 torsion angles that minimise this protein-specific potential using only limited sampling. We also
92 show that whole chains can be optimised together, avoiding the need for segmenting long proteins
93 into hypothesised domains which are modelled independently.

94 The central component of AlphaFold is a convolutional neural network which is trained
95 on PDB structures to predict the distances d_{ij} between the C_β atoms of pairs, ij , of a protein's
96 residues. Based on a representation of the protein's amino acid sequence, \mathcal{S} , and features derived
97 from the sequence's MSA, the network, similar in structure to those used for image recognition
98 tasks³¹, predicts a discrete probability distribution $P(d_{ij} | \mathcal{S}, \text{MSA}(\mathcal{S}))$ for every ij pair in a
99 64×64 residue region, as shown in Fig. 2b. The full set of distance distribution predictions
100 is constructed by averaging predictions for overlapping regions and is termed a *distogram* (from
101 distance histogram). Figure 3 shows an example histogram prediction for one CASP protein,
102 T0955. The modes of the distribution (Fig. 3c) can be seen to closely match the true distances
103 (Fig. 3b). Example distributions for all distances to one residue (29) are shown in Fig. 3c. Further
104 analysis of how the network predicts the distances is shown in Methods Figure 14.

105 In order to realise structures that conform to the distance predictions, we construct a smooth
106 potential V_{distance} by fitting a spline to the negative log probabilities, and summing across all the
107 residue pairs. We parameterise protein structures by the backbone torsion angles (ϕ, ψ) of all
108 residues and build a differentiable model of protein geometry $\mathbf{x} = G(\phi, \psi)$ to compute the C_β
109 coordinates, \mathbf{x} , and thus the inter-residue distances, $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|$, for each structure, and
110 express V_{distance} as a function of ϕ and ψ . For a protein with L residues, this potential accumulates
111 L^2 terms from marginal distribution predictions. To correct for the over-representation of the
112 prior we subtract a *reference distribution*³² from the distance potential in the log domain. The
113 reference distribution models the distance distributions $P(d_{ij} | \text{length})$ independent of the protein
114 sequence and is computed by training a small version of the distance prediction neural network on
115 the same structures, without sequence or MSA input features. A separate output head of the contact
116 prediction network is trained to predict discrete probability distributions of backbone torsion angles
117 $P(\phi_i, \psi_i | \mathcal{S}, \text{MSA}(\mathcal{S}))$. After fitting a von Mises distribution, this is used to add a smooth torsion
118 modelling term $V_{\text{torsion}} = -\sum \log p_{\text{vonMises}}(\phi_i, \psi_i | \mathcal{S}, \text{MSA}(\mathcal{S}))$ to the potential. Finally, to
119 prevent steric clashes, we add Rosetta's $V_{\text{score2_smooth}}$ ¹⁰ to the potential, as this incorporates a van
120 der Waals term. We used multiplicative weights for each of the three terms in the potential, but no
121 weighting noticeably outperformed equal weighting.