

AlphaFold: Improved protein structure prediction using potentials from deep learning

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

¹DeepMind, London, UK

²The Francis Crick Institute, London, UK

³University College London, London, UK

*These authors contributed equally to this work.

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

Proteins are at the core of most biological processes. Since the function of a protein is dependent on its structure, understanding protein structure has been a grand challenge in biology 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.

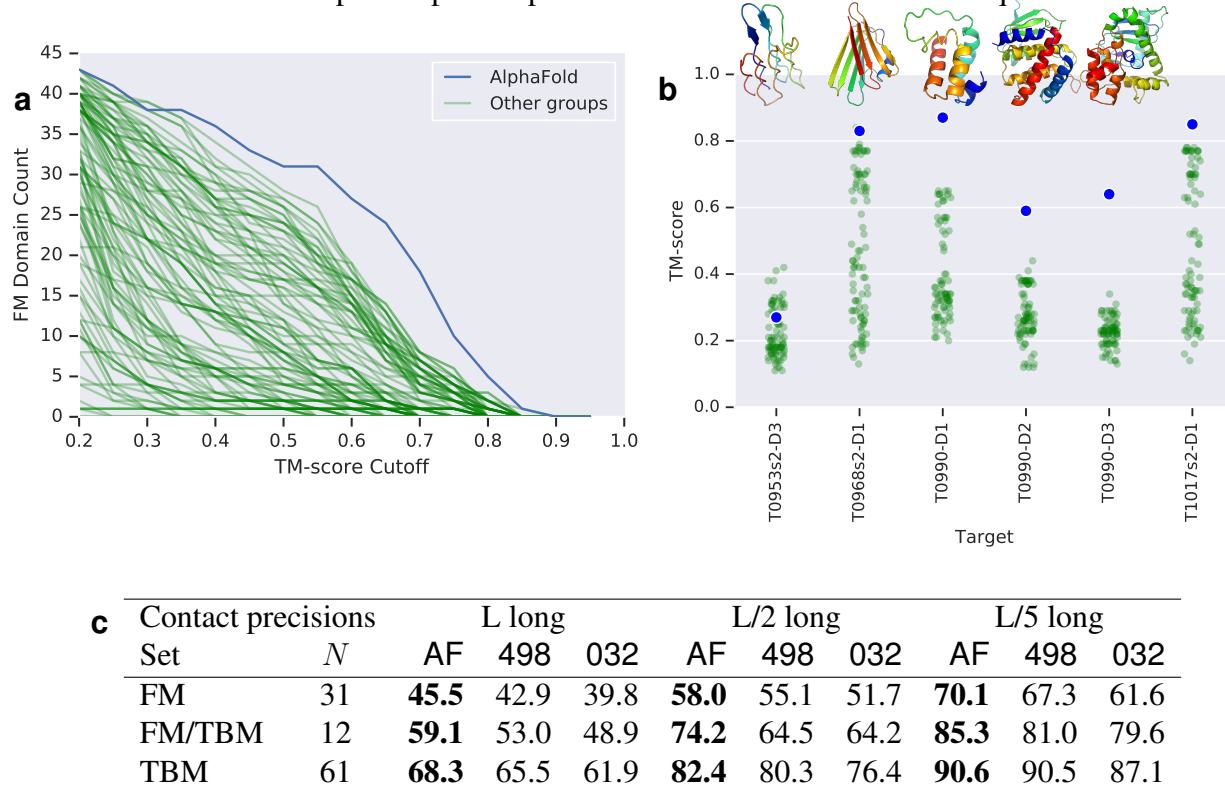


Fig. 1 | AlphaFold’s performance in the CASP13 assessment. (a) Number of free modelling (FM + FM/TBM) domains predicted to a given TM-score threshold for AlphaFold and the other 97 groups. (b) For the six new folds identified by the CASP13 assessors, AlphaFold’s TM-score compared with the other groups, 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

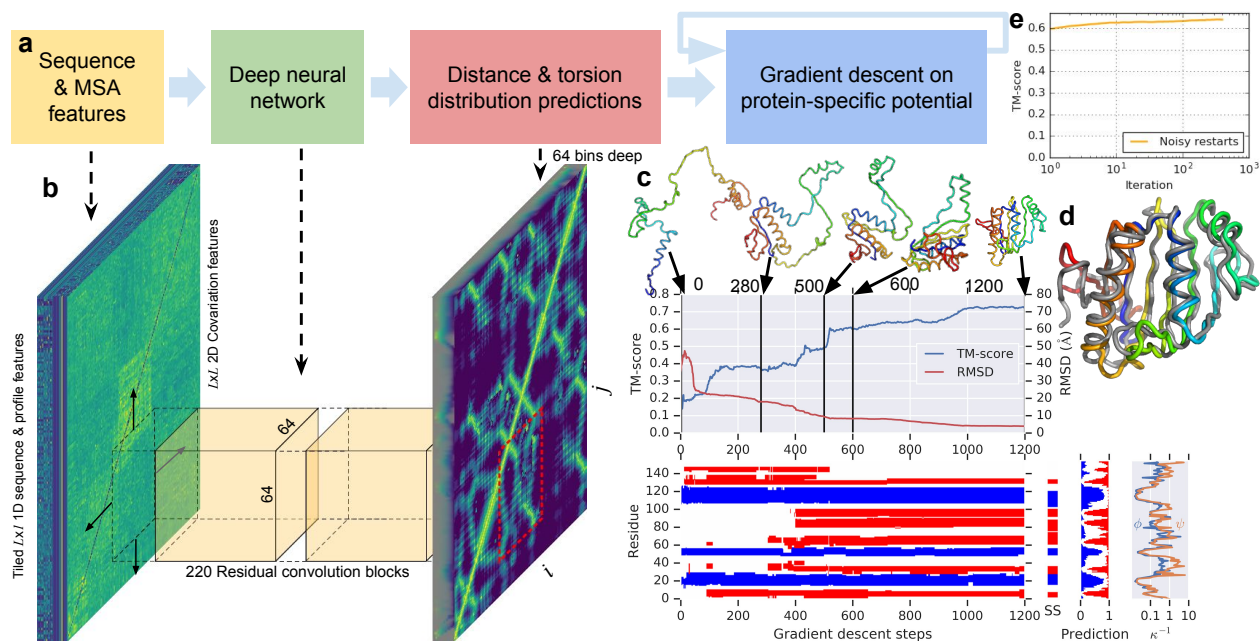


Fig. 2 | The folding process illustrated for CASP13 target T0986s2. (Length $L = 155$) (a) Steps of structure prediction. (b) The neural network predicts the entire $L \times L$ distogram based on MSA features, accumulating separate predictions for 64×64 -residue regions. (c) One iteration of gradient 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 $P(d_{ij} \mid \mathcal{S}, \text{MSA}(\mathcal{S}))$ 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 distogram 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 $P(d_{ij} \mid \text{length})$ 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 $P(\phi_i, \psi_i \mid \mathcal{S}, \text{MSA}(\mathcal{S}))$. After fitting a von Mises distribution, this is used to add a smooth torsion
118 modelling term $V_{\text{torsion}} = -\sum \log p_{\text{vonMises}}(\phi_i, \psi_i \mid \mathcal{S}, \text{MSA}(\mathcal{S}))$ 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.

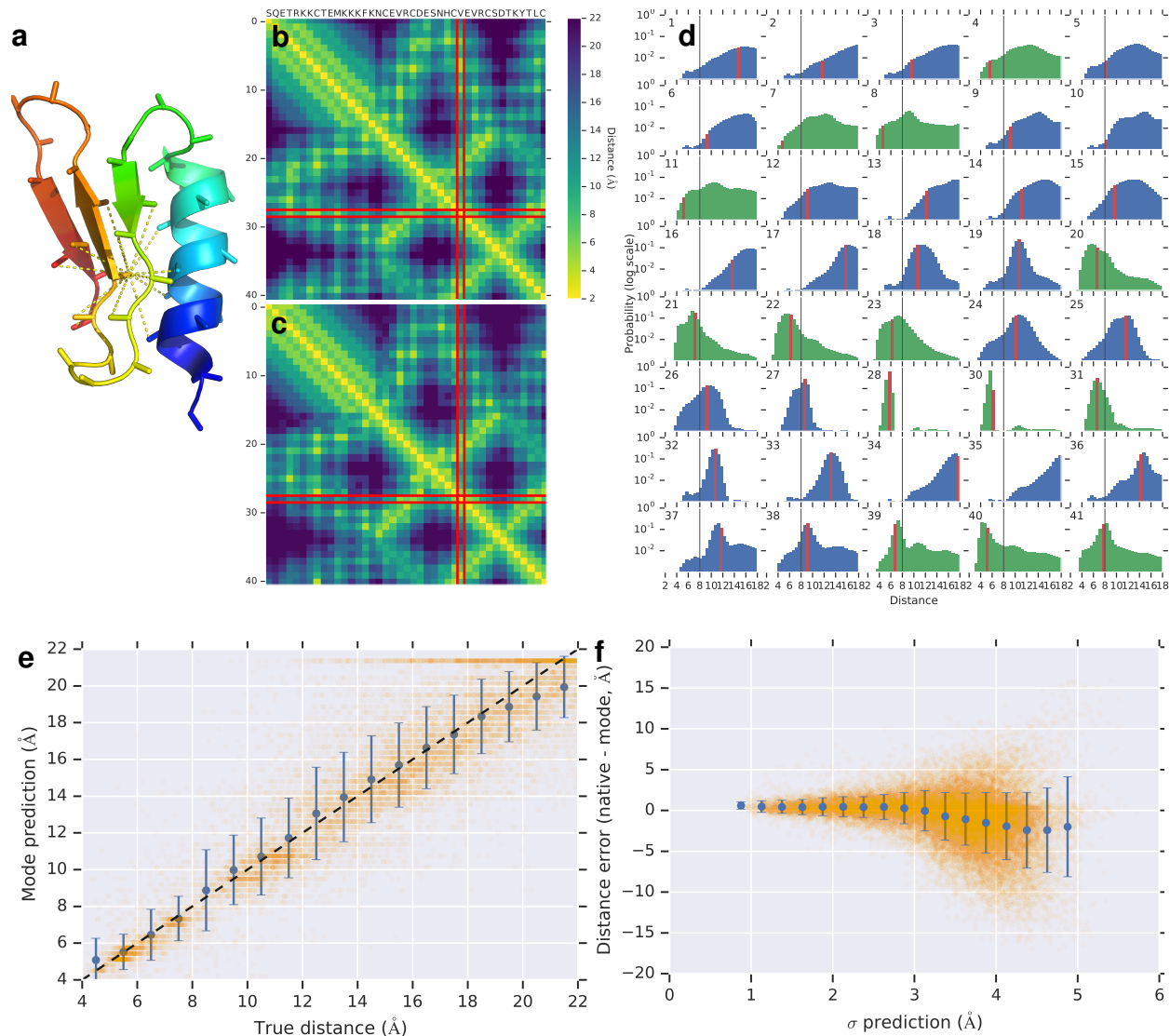


Fig. 3 | Predicted distance distributions compared with true distances. Above, for CASP target T0955 ($L = 41$): (a) Native structure showing distances under 8 \AA from C_β of residue 29. (b) Native inter-residue distances and (c) the mode of the distance predictions, highlighting residue 29. (d) The predicted probability distributions for distances of residue 29 to all other residues. The bin corresponding to the native distance is highlighted in red, 8 \AA drawn in black. True contacts' distributions are plotted in green, non-contacts in blue. Below, for CASP target T0990 ($L = 552$): (e) the mode of the predicted distance plotted against the true distance for all residue pairs with distances $\leq 22 \text{ \AA}$, excluding distributions with standard deviation $> 3.5 \text{ \AA}$. The blue error bars show mean and standard deviation calculated for 1 \AA bins. (f) The error of the mode distance prediction vs the standard deviation of the distance distributions, excluding pairs with native distances $> 22 \text{ \AA}$. Mean and standard deviations are shown for 0.25 \AA bins. The distogram is shown in Figure 7 in Methods.

122 Since all the terms in the combined potential $V_{\text{total}}(\phi, \psi)$ are differentiable functions of
123 (ϕ, ψ) , it can be optimised with respect to these variables by gradient descent. Here we use
124 L-BFGS³³. Structures are initialised by sampling torsion values from $P(\phi_i, \psi_i \mid \mathcal{S}, \text{MSA}(\mathcal{S}))$.
125 Figure 2c illustrates a single gradient descent trajectory minimising the potential, showing how this
126 greedy optimisation process leads to increasing accuracy and large-scale conformation changes.
127 Secondary structure is partly set by the initialisation, since some areas of secondary structure are
128 predicted accurately, leading to low-variance torsion angle distributions. Overall accuracy (TM-
129 score) improves quickly and after a few hundred steps of gradient descent has converged to a local
130 optimum.

131 We repeat the optimisation from sampled initialisations, leading to a pool of low potential
132 structures from which further structure initialisations are sampled, with added backbone torsion
133 noise (‘noisy restarts’), leading to more structures to be added to the pool. After only a few
134 hundred cycles the optimisation converges and the lowest potential structure is chosen as the best
135 candidate structure. Figure 2e shows the progress in the accuracy of the best-scoring structures over
136 multiple restarts of the gradient descent process, showing that after a few iterations the optimisation
137 has converged. Noisy restarts enable slightly higher TM-score structures to be found than when
138 continuing to sample from the predicted torsion distributions (average of 0.641 vs 0.636 on our test
139 set).

140 A key component of AlphaFold’s overall accuracy is that accurate distance predictions con-
141 vey more information about structure than contact predictions. Figure 3e shows that the predictions
142 of distance correlate well with the true distance. It can be seen from Figure 3f that the network
143 is also modelling the uncertainty in its predictions. When the standard deviation of the predicted
144 distribution is low, the predictions are more accurate. This is also evident in the predicted distri-
145 butions of Figure 3d, where more confident predictions of the distance distribution (higher peak
146 and lower standard deviation of the distribution) tend to be more accurate, with the true distance
147 close to the peak. Broader, less-confidently-predicted distributions still assign probability to the
148 correct value even when it is not close to the peak. The high accuracy of the distance predictions
149 and consequently the contact predictions (Table 1c) comes from a combination of factors in the de-
150 sign of the neural network and its training, including predicting distances instead of contacts, data
151 augmentation, feature representation, auxiliary losses, cropping and data curation. (See Methods
152 section.)

153 Figure 4a shows that the distogram accuracy (measured by distogram IDDT₁₂, defined in
154 Methods) correlates well with the TM-score of the final realised structures. Figure 4b shows the
155 effect of changing the construction of the potential. Removing the distance potential entirely gives
156 a TM-score of 0.266. Reducing the resolution of the distogram representation below 6 bins by av-
157 eraging adjacent bins causes the TM-score to degrade. Removing the torsion potential, reference
158 correction or $V_{\text{score2_smooth}}$ degrade the accuracy only slightly. A final ‘relaxation’ (side-chain pack-
159 ing interleaved with gradient descent) with Rosetta¹⁰, using a combination of the Talaris2014
160 potential and a spline fit of our reference-corrected distance potential adds side-chain atom coor-

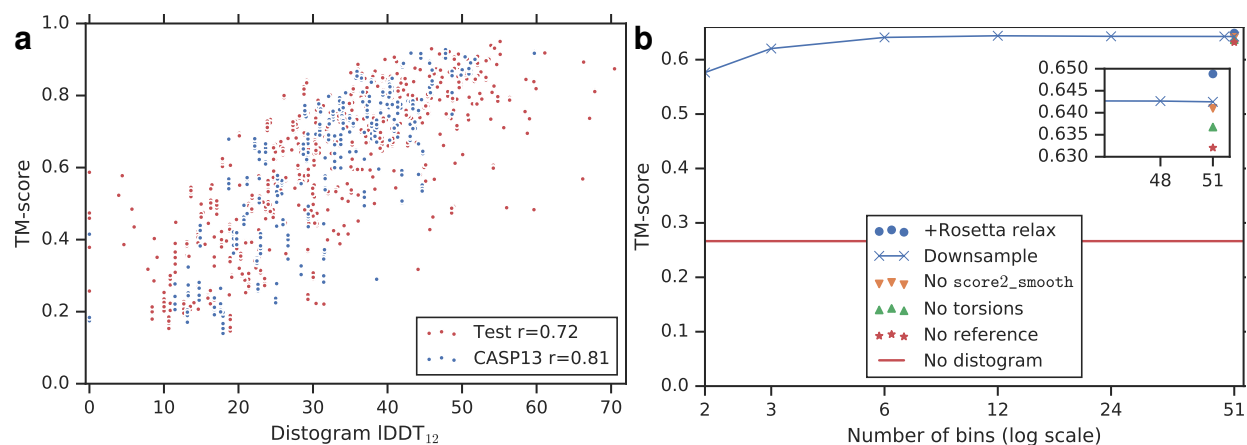


Fig. 4 | TM-scores vs the accuracy of the distogram, and the TM scores' dependency on different components of the potential. (a) TM-score vs distogram IDDT₁₂ with Pearson's correlation coefficients, for both CASP13 ($n = 108$) and test ($n = 377$) datasets. (b) Average TM-score over the test set ($n = 377$) vs number of histogram bins used when downsampling the distogram, compared with removing different components of the potential, or adding Rosetta relaxation.

161 dinates, and yields a small average improvement of 0.007 TM-score.

162 We have shown that a carefully designed deep-learning system can provide accurate predic-
 163 tions of inter-residue distances and be used to construct a protein-specific potential which repre-
 164 sents protein structure. Furthermore we have shown that this potential can be simply optimised
 165 with gradient descent to achieve accurate structure predictions. While free modelling predictions
 166 only rarely approach the accuracy of experimental structures, the CASP13 assessment shows that
 167 the AlphaFold system achieves unprecedented free modelling accuracy and that this free modelling
 168 method can match the performance of template modelling approaches without using templates and
 169 is starting to reach the accuracy needed for biological understanding (see Methods). We hope that
 170 the methods we have described can be developed further and applied to benefit all areas of protein
 171 science with more accurate predictions for sequences of unknown structure.

172 References

- 173 1. Dill, K., Ozkan, S. B., Shell, M. & Weikl, T. The protein folding problem. *Annu. Rev. Biophys.*
 174 **37**, 289–316 (2008).
- 175 2. Dill, K. & MacCallum, J. The protein-folding problem, 50 years on. *Science* **338**, 1042–1046
 176 (2012).
- 177 3. Schaarschmidt, J., Monastyrskyy, B., Kryshchak, A. & Bonvin, A. M. Assessment of
 178 contact predictions in CASP12: Co-evolution and deep learning coming of age. *Proteins* **86**,
 179 51–66 (2018).

- 180 4. Kirkwood, J. Statistical mechanics of fluid mixtures. *J. Chem. Phys.* **3**, 300–313 (1935).
- 181 5. Moulton, J., Fidelis, K., Kryshtafovych, A., Schwede, T. & Tramontano, A. Critical assess-
182 ment of methods of protein structure prediction (CASP)—Round XII. *Proteins: Structure,*
183 *Function, and Bioinformatics* **86**, 7–15 (2018).
- 184 6. Zhang, Y. & Skolnick, J. Scoring function for automated assessment of protein structure
185 template quality. *Proteins* **57**, 702–710 (2004).
- 186 7. Zhang, Y. Protein structure prediction: when is it useful? *Current opinion in structural biol-*
187 *ogy* **19**, 145–155 (2009).
- 188 8. Xu, J. *Protein structure modeling by predicted distance instead of contacts in CASP13 Ab-*
189 *stracts* Dec. 1, 2018 (2018), 146–7.
- 190 9. Zhang, C., Li, Y., Yu, D. & Zhang, Y. *Contact map prediction by deep residual fully convo-*
191 *lutional neural network with only evolutionary coupling features derived from deep multiple*
192 *sequence alignment in CASP13 Abstracts* Dec. 1, 2018 (2018), 181–2.
- 193 10. Das, R. & Baker, D. Macromolecular modeling with Rosetta. *Annu. Rev. Biochem.* **77**, 363–
194 382 (2008).
- 195 11. Jones, D. T. Predicting novel protein folds by using FRAGFOLD. *Proteins* **45**, 127–132
196 (2001).
- 197 12. Zhang, C., Mortuza, S., He, B., Wang, Y. & Zhang, Y. Template-based and free modeling
198 of I-TASSER and QUARK pipelines using predicted contact maps in CASP12. *Proteins* **86**,
199 136–151 (2018).
- 200 13. Kirkpatrick, S., Gelatt, C. & Vecchi, M. Optimization by simulated annealing. *Science* **220**,
201 671–680 (1983).
- 202 14. Gilliland, G. *et al.* The Protein Data Bank. *Nucleic Acids Res.* **28**, 235–242 (2000).
- 203 15. Altschuh, D., Lesk, A., Bloomer, A. & Klug, A. Correlation of co-ordinated amino acid
204 substitutions with function in viruses related to tobacco mosaic virus. *J. Mol. Biol.* **193**, 693–
205 707 (1987).
- 206 16. Ovchinnikov, S., Kamisetty, H. & Baker, D. Robust and accurate prediction of residue–
207 residue interactions across protein interfaces using evolutionary information. *Elife* **3**, e02030
208 (2014).
- 209 17. Seemayer, S., Gruber, M. & Söding, J. CCMpred—fast and precise prediction of protein
210 residue–residue contacts from correlated mutations. *Bioinformatics* **30**, 3128–3130 (2014).
- 211 18. Morcos, F. *et al.* Direct-coupling analysis of residue coevolution captures native contacts
212 across many protein families. *PNAS* **108**, E1293–E1301. ISSN: 0027-8424 (2011).
- 213 19. Jones, D. T., Buchan, D. W., Cozzetto, D. & Pontil, M. PSICOV: Precise structural contact
214 prediction using sparse inverse covariance estimation on large multiple sequence alignments.
215 *Bioinformatics* **28**, 184–190 (2011).

- 216 20. Skwark, M., Raimondi, D., Michel, M. & Elofsson, A. Improved Contact Predictions Using
217 the Recognition of Protein Like Contact Patterns. *PLoS Computational Biology* **10**, 1–14
218 (2014).
- 219 21. Jones, D., Singh, T., Kosciolok, T. & Tetchner, S. MetaPSICOV: Combining coevolution
220 methods for accurate prediction of contacts and long range hydrogen bonding in proteins.
221 *Bioinformatics* **31**, 999–1006 (2015).
- 222 22. Wang, S., Sun, S., Li, Z., Zhang, R. & Xu, J. Accurate De Novo Prediction of Protein Contact
223 Map by Ultra-Deep Learning Model. *PLoS Comput. Biol.* **13**, 999–1006 (2017).
- 224 23. Jones, D. T. & Kandathil, S. M. High precision in protein contact prediction using fully
225 convolutional neural networks and minimal sequence features. *Bioinformatics* **1**, 8 (2018).
- 226 24. Ovchinnikov, S. *et al.* Improved de novo structure prediction in CASP 11 by incorporating
227 coevolution information into Rosetta. *Proteins* **84**, 67–75 (2016).
- 228 25. Aszódi, A. & Taylor, W. R. Estimating polypeptide α -carbon distances from multiple se-
229 quence alignments. *J. Math. Chem.* **17**, 167–184 (1995).
- 230 26. Zhao, F. & Xu, J. A position-specific distance-dependent statistical potential for protein struc-
231 ture and functional study. *Structure* **20**, 1118–1126 (2012).
- 232 27. Aszodi, A., Gradwell, M. & Taylor, W. Global fold determination from a small number of
233 distance restraints. *J. Mol. Biol.* **251**, 308–326 (1995).
- 234 28. Kandathil, S., Greener, J. & Jones, D. *DMPfold: a new deep learning-based method for*
235 *protein tertiary structure prediction and model refinement in CASP13 Abstracts* Dec. 1, 2018
236 (2018), 84–5.
- 237 29. Xu, J. Distance-based Protein Folding Powered by Deep Learning. *arXiv preprint arXiv:1811.03481*
238 (2018).
- 239 30. Konagurthu, A. S., Lesk, A. M. & Allison, L. Minimum message length inference of sec-
240 ondary structure from protein coordinate data. *Bioinformatics* **28**, i97–i105 (2012).
- 241 31. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition. *arXiv*
242 *preprint arXiv:1512.03385* **abs/1512.03385** (2015).
- 243 32. Simons, K., Kooperberg, C., Huang, E. & Baker, D. Assembly of Protein Tertiary Structures
244 from Fragments with Similar Local Sequences using Simulated Annealing and Bayesian
245 Scoring Functions. *J. Mol. Biol.* **268**, 209–225 (1997).
- 246 33. Liu, D. C. & Nocedal, J. On the limited memory BFGS method for large scale optimization.
247 *Math. Program.* **45**, 503–528 (1989).

248 **1 Methods**

249 Figure 5 shows the steps involved in MSA construction, feature extraction, distance prediction,
250 potential construction and structure realisation.

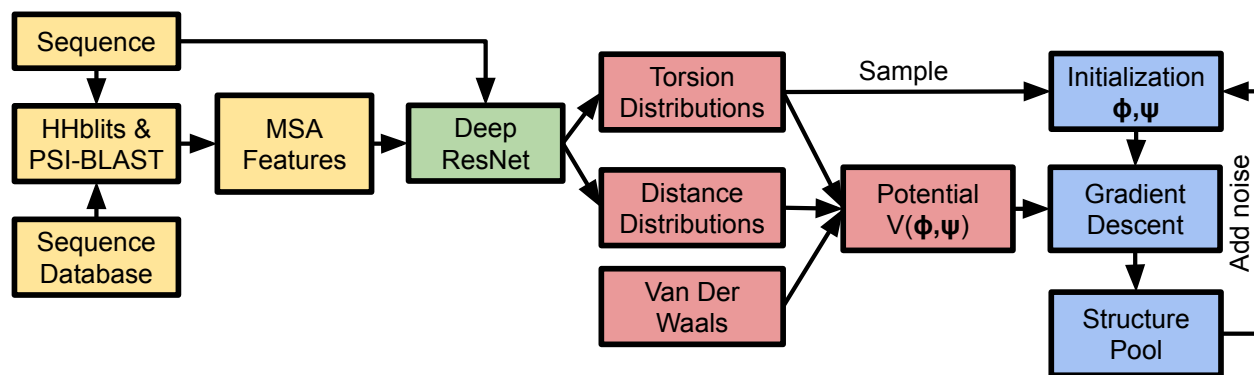


Fig. 5 | A schematic of the folding system. Feature extraction stages are shown in yellow, structure-prediction neural network in green, potential construction in red and structure realisation in blue.

251 **Data** Our models are trained on structures extracted from the Protein Data Bank¹. We extract
 252 non-redundant domains by utilising the CATH² 35% sequence similarity cluster representatives.
 253 This gives 31 247 domains, which are split into train, and test sets (29 427 and 1 820 proteins
 254 respectively) keeping all domains from the same homologous superfamily (H-level in the CATH
 255 classification) in the same partition. The CATH superfamilies of FM domains from CASP11 and
 256 CASP12 were also excluded from the training set. From the test set, we take a single domain per
 257 homologous superfamily to create the 377 domain subset used for the results presented here. We
 258 note that accuracies for this set are higher than for the CASP13 test domains.

259 CASP13 submission results are drawn from the CASP13 results pages with additional re-
 260 sults shown for the CASP13 dataset for “all groups” chains, scored on CASP13 PDB files, by
 261 CASP domain definitions. Contact prediction accuracies are recomputed from the group 032 and
 262 498 submissions (as RR files), compared with the distogram predictions used by AlphaFold for
 263 CASP13 submissions. Contact prediction probabilities are obtained from the distograms by sum-
 264 ming the probability mass in each distribution below 8 Å.

265 For each training sequence, we search for and align similar protein sequences in the Uni-
 266 clust30³ dataset with HHblits⁴ and use the returned MSA to generate *profile* features with the
 267 position-specific substitution probabilities for each residue as well as covariation features — the
 268 parameters of a regularised pseudolikelihood-trained Potts model similar to CCMPred⁵. CCMPred
 269 uses the Frobenius norm of the parameters, but we feed both this norm (1 feature) as well as the
 270 raw parameters (484 features) into the network for each residue pair ij . In addition we provide
 271 the network with features explicitly representing gaps and deletions in the MSA. To make the net-
 272 work better able to make predictions for shallow MSAs, and as a form of data augmentation, we
 273 take a sample of half the sequences from the the HHblits MSA before computing the MSA-based
 274 features. Our training set contains 10 such samples for each domain. We extract additional profile
 275 features using PSI-BLAST⁶.

276 The distance prediction neural network was trained with the following input features (with
277 number of features).

- 278 • Number of HHblits alignments (1D scalar)
- 279 • Sequence-length features: 1-hot amino acid type (21D), Profiles: PSI-BLAST (21D), HH-
280 blits profile (22D), non-gapped profile (21D), HHblits bias, HMM profile (30D) Potts model
281 bias (22D); Deletion probability (1D); residue index (integer index of residue number, con-
282 secutive except for multi-segment domains, encoded as 5 least-significant bits and a scalar);
- 283 • Sequence-length-squared features: Potts model parameters (484D, fitted with 500 iterations
284 of gradient descent using Nesterov momentum 0.99, without sequence reweighting); Frobe-
285 nius norm (1D); Gap matrix (1D)

286 **Distogram prediction** The inter-residue distances are predicted by a deep neural network. The
287 architecture is a deep two-dimensional dilated convolutional residual network. Xu *et al.*⁷ used
288 a two-dimensional residual network preceded by one-dimensional embedding layers for contact
289 prediction. Our network is two-dimensional throughout and uses 220 residual blocks⁸ with dilated
290 convolutions⁹. Each residual block, illustrated in Figure 6 consists of a sequence of neural network
291 layers¹⁰, interleaving three batchnorm layers; two 1×1 projection layers; a 3×3 dilated convolution
292 layer and ELU¹¹ nonlinearities. Successive layers cycle through dilations of 1, 2, 4, 8 pixels to
293 allow propagation of information quickly across the cropped region. At the final layer, a position-
294 specific bias was used, so the biases were indexed by residue-offset (capped at 32) and bin number.

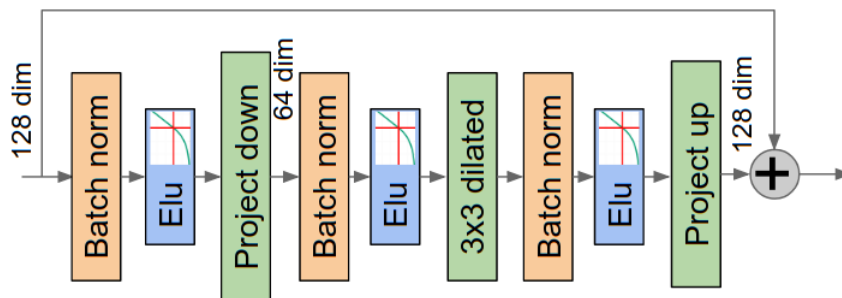


Fig. 6 | The layers used in one block of the deep residual convolutional network. The dilated convolution is applied on reduced-dimensional features. The output of the block is added to the representation from the previous layer. The residual network's bypass connections allow gradients to pass back through the network undiminished, permitting the training of very deep networks.

295

296 The network is trained with stochastic gradient descent using a cross-entropy loss. The target
297 is a quantisation of the distance between the residues' C_β atoms (C_α for glycine). We quantise the
298 range 2–22 Å into 64 equal bins. The input to the network consists of a two-dimensional array of
299 features where each i, j feature is the concatenation of the 1-dimensional features for both i and j

300 as well as the two-dimensional features for i, j .

301 Individual training runs were cross-validated with early stopping using 27 CASP11 FM do-
302 mains as a validation set. Models were selected by cross-validation on 27 CASP12 FM domains.

303 **Neural network hyperparameters**

- 304 • 7×4 Blocks with 256 channels, cycling through dilations 1, 2, 4, 8
- 305 • 48×4 Blocks with 128 channels, cycling through dilations 1, 2, 4, 8
- 306 • Optimisation: Synchronized stochastic gradient descent
- 307 • Batch size: batch of 4 crops on each of 8 GPU workers
- 308 • 0.85 Dropout keep probability
- 309 • Nonlinearity: ELU
- 310 • Learning rate 0.06
- 311 • Auxilliary loss weights: Secondary structure: 0.005; Accessible surface area: 0.001. These
312 auxilliary losses were cut by a factor 10 after 100 000 steps.
- 313 • Learning rate decayed by 50% at 150 000, 200 000, 250 000 and 350 000 steps.
- 314 • Training time: about 5 days for 600 000 steps

315 To constrain memory usage and avoid overfitting, the network is always trained on 64×64
316 regions of the distance matrix, that is the pairwise distances between 64 consecutive residues and
317 another group of 64 consecutive residues. For each training domain, the entire distance matrix is
318 split into non-overlapping 64×64 crops. By training off-diagonal crops, the interaction between
319 residues further apart than 64 residues can be modelled. Each crop consists of the distance matrix
320 which represents the juxtaposition of two 64-residue fragments. Jones and Kandathil¹² have shown
321 that contact prediction needs only a limited context window. We note that the distance predictions
322 close to the diagonal $i = j$, encode predictions of the local structure of the protein, and for any
323 cropped region the distances are governed by the local structure of the two fragments represented
324 by the i and j ranges of the crop. Augmenting the inputs with the on-diagonal 2D input features
325 that correspond to both the i and j ranges provides additional information to predict the structure of
326 each fragment and thus distances between them. It can be seen that if the fragment structures can
327 be well predicted (for instance if they are confidently predicted as helices or sheets) then prediction
328 of a single contact between the fragments will strongly constrain the distances between all other
329 pairs

330 Randomising the offset of the crops each time a domain is used in training leads to a form
331 of data augmentation where a single protein can generate many thousands of different training
332 examples. This is further enhanced by adding noise to the atom coordinates, proportional to the
333 ground truth resolution leading to variation in the target distances. Data augmentation (MSA
334 subsampling and coordinate noise), together with dropout, prevents the network from overfitting
335 to the training data.

336 To predict the distance distribution for all $L \times L$ residue pairs, many 64×64 crops are com-
 337 bined. To avoid edge effects, several such tilings are produced with different offsets and averaged
 338 together, with a heavier weighting for the predictions near the centre of the crop. To improve
 339 accuracy further, predictions from an ensemble of four separate models, trained independently
 340 with slightly different hyperparameters, are averaged together. Figure 7 shows an example of true
 341 distances (a) and the mode of the histogram prediction (b) for a three-domain CASP13 target.

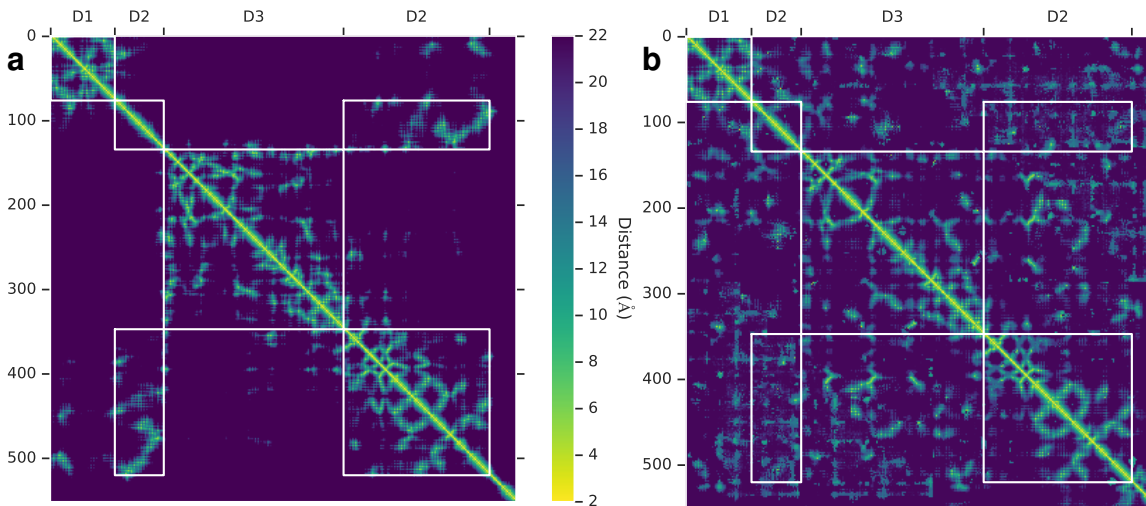
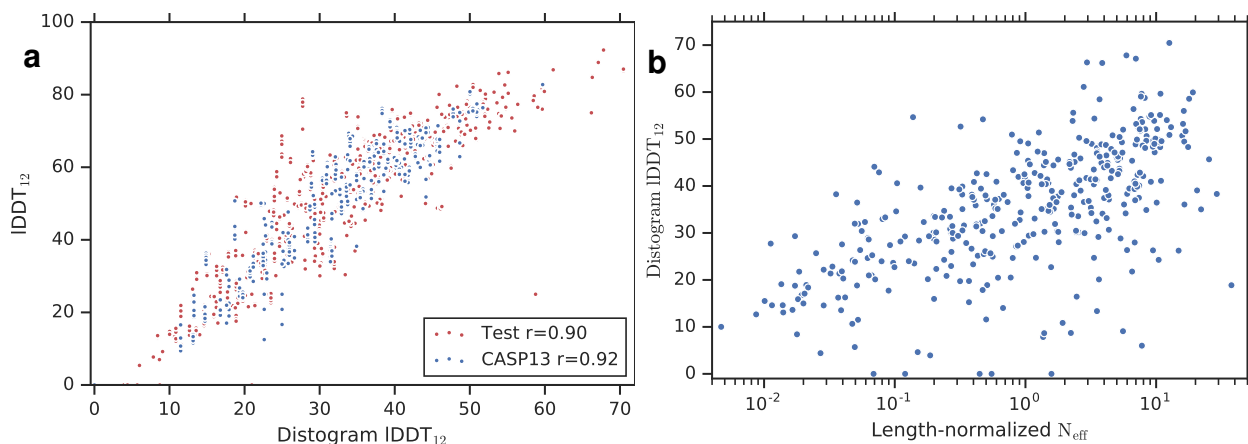


Fig. 7 | True distances (a) and modes of the predicted histogram (b) for CASP13 target T0990. CASP divides this chain into 3 domains as shown (D3 is inserted in D2) for which there are 39, 36 and 42 HHblits alignments respectively (from the CASP website).

342 Since the network has a rich representation capable of incorporating both profile and covari-
 343 ation features of the MSA, we argue that the network can be used to predict secondary structure
 344 directly. By mean- and max-pooling the 2D activations of the penultimate layer of the network
 345 separately in both i and j , we add an additional 1-dimensional output head to the network which
 346 predicts 8-class secondary structure labels as computed by DSSP¹³ for each residue in j and i .
 347 The resulting Q3 (distinguishing the three helix / sheet / coil classes) predictions' accuracy is 84%
 348 which is comparable to the state-of-the-art¹⁴. The relative accessible surface area (ASA) of each
 349 residue can also be predicted.

350 The 1-dimensional pooled activations are also used to predict the marginal Ramachandran
 351 distributions: $P(\phi_i, \psi_i \mid \mathcal{S}, \text{MSA}(\mathcal{S}))$, independently for each residue, as a discrete probability
 352 distribution quantised to 10° (1296 bins). In practice during CASP13 we used histograms from a
 353 network that was trained to predict histograms, secondary structure and ASA with torsions from
 354 a second, similar network trained to predict histograms, secondary structure, ASA and torsions,
 355 since the former had been more thoroughly validated.

356 Figure 8b shows that an important factor in the accuracy of the histograms (as has previ-
 357 ously been found with contact prediction systems) is N_{eff} , the effective number of sequences in the



Potential	Bins	TM-score	GDT_TS	IDDT	RMSD (Å)
Full + relax	51/64	0.649	65.8	54.2	5.94
Full	51/64	0.642	65.0	53.9	5.91
W/o reference	51/64	0.632	64.3	50.0	6.64
W/o score2_smooth	51/64	0.641	64.8	53.7	5.93
W/o torsions	51/64	0.637	64.3	53.6	6.04
W/o distogram	51/64	0.266	29.1	19.1	14.88
Full	48/64	0.643	65.0	54.1	5.90
Full	24/32	0.643	65.0	53.8	5.89
Full	12/16	0.644	65.1	53.9	5.85
Full	6/8	0.641	64.6	53.7	5.94
Full	3/4	0.620	62.4	52.8	6.22
Full	2/3	0.576	58.2	49.3	8.38

Fig. 8 | Analysis of structure accuracies. (a) IDDT vs distogram IDDT₁₂ (Defined below under ‘Accuracy’). The distogram accuracy predicts the realised structure’s IDDT (as well as TM-score as shown in Fig. 4a) well for both CASP13 ($n = 108$) and test ($n = 377$) datasets. Shown with Pearson’s correlation coefficients. (b) DLDDT₁₂ against effective number of sequences in the MSA (N_{eff}) normalised by sequence length ($n = 377$). The number of effective sequences correlates with this measure of distogram accuracy ($r = 0.634$). (c) Structure accuracy measures, computed on the test set, for gradient descent optimisation of different forms of the potential. Above: Removing terms in the potential, also showing the effect of following optimisation with Rosetta relax. Bins shows the number of bins fitted by the spline before extrapolation and the number in the full distribution. In CASP13 splines were fitted to the first 51 of 64 bins. Below, reducing the resolution of the distogram distributions. The original 64-bin distogram predictions are repeatedly downsampled by a factor 2 by summing adjacent bins, in each case with constant extrapolation beyond 18 Å (the last $\frac{1}{4}$ of the bins). The final row’s two-level potential, designed to compare to contact predictions, is constructed by summing the probability mass below 8 Å and between 8–14 Å, with constant extrapolation beyond 14 Å. The TM-scores in this table are plotted in Figure 4 (b) Accuracy measures, computed on the test set ($n = 377$), for gradient descent optimisation with differently constructed potentials. The TM-scores in this table are plotted in Figure 4b.

358 MSA¹⁵. This is the number of sequences found in the MSA, discounting redundancy at the 62%
359 sequence identity level, which we then divide by the number of residues in the target, and is an
360 indication of the amount of covariation information in the MSA.

361 **Distance potential** The histogram probabilities are estimated for discrete distance bins, so to con-
362 struct a differentiable potential the distribution is interpolated with a cubic spline. Because the final
363 bin accumulates probability mass from all distances beyond 22 Å, and since greater distances are
364 harder to predict accurately, the potential is only fit up to 18 Å (determined by cross-validation),
365 with a constant extrapolation thereafter.

366 To predict a reference distribution, a similar model is trained on the same dataset. The
367 reference distribution is not conditioned on the sequence, but to account for the atoms between
368 which we are predicting distances, we do provide a feature $\delta_{\alpha\beta}$ to indicate if the residue is glycine
369 (C_α atom) or not (C_β) and the overall length of the protein.

370 A distance potential is created from the negative log likelihood of the distances, summed
371 over all pairs of residues i, j .

$$V_{\text{distance}}(\mathbf{x}) = - \sum_{i,j, i \neq j} \log P(d_{ij} | \mathcal{S}, \text{MSA}(\mathcal{S})) \quad (1)$$

372 With a reference state this becomes the log likelihood ratio of the distances under the full condi-
373 tional model and under the background model:

$$V_{\text{distance}}(\mathbf{x}) = - \sum_{i,j, i \neq j} \log P(d_{ij} | \mathcal{S}, \text{MSA}(\mathcal{S})) - \log P(d_{ij} | \text{length}, \delta_{\alpha\beta}) \quad (2)$$

374 Torsions are modelled as a negative log likelihood under the predicted torsion distributions.
375 Since we have marginal distribution predictions, each of which can be multimodal, it can be dif-
376 ficult to jointly optimise the torsions. To unify all the probability mass, at the cost of modelling
377 fidelity of multimodal distributions, we fit a unimodal von Mises distribution to the marginal pre-
378 dictions. The potential is summed over all residues i .

$$V_{\text{torsion}}(\phi, \psi) = - \sum_i \log p_{\text{vonMises}}(\phi_i, \psi_i | \mathcal{S}, \text{MSA}(\mathcal{S})) \quad (3)$$

379 Finally, to prevent steric clashes, a van der Waals term is introduced through the use of
380 Rosetta's $V_{\text{score2.smooth}}$.

381 **Structure realisation by gradient descent** To realise structures which minimise the constructed
382 potential, we create a differentiable model of ideal protein backbone geometry, giving backbone

383 atom coordinates as a function of the torsion angles (ϕ, ψ) : $\mathbf{x} = G(\phi, \psi)$. The complete potential
384 to be minimised is then*:

$$V_{\text{total}}(\phi, \psi) = V_{\text{distance}}(G(\phi, \psi)) + V_{\text{torsion}}(\phi, \psi) + V_{\text{score2.smooth}}(G(\phi, \psi)). \quad (4)$$

385 Since every term in V_{total} is differentiable with respect to the torsion angles, given an initial set of
386 torsions ϕ, ψ which can be sampled from the predicted torsion marginals, we can minimise V_{total}
387 using a gradient descent algorithm, such as L-BFGS¹⁶. The optimised structure is dependent on the
388 initial conditions, so we repeat the optimisation multiple times with different initialisations. A pool
389 of the 20 lowest-potential structures is maintained and once full, we initialise 90% of trajectories
390 from those with 30° noise added to the backbone torsions (the remaining 10% being sampled from
391 the predicted torsion distributions). In CASP13 we made 5000 optimisation runs for each chain.
392 Figure 2 shows the change in TM-score against the number of restarts. Since longer chains take
longer to optimise, this work load was balanced across $(50 + L)/2$ parallel workers. Figure 9 shows

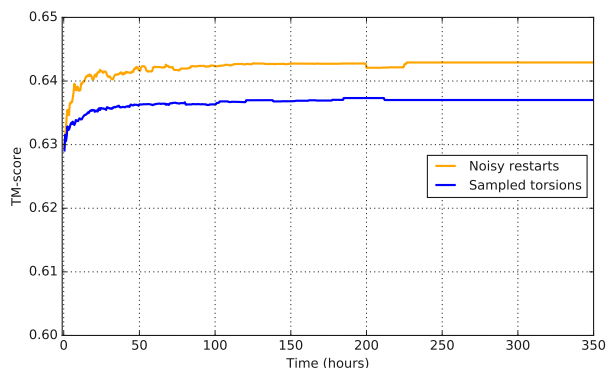


Fig. 9 | TM-score vs per-target computation time computed as an average over the test set ($n = 377$). Full optimisation with noisy restarts (orange) is compared with initialisation from sampled torsions (blue). Computation is measured as the product of the number of (CPU-based) machines and time elapsed and can be largely parallelised. Longer targets take longer to optimise.

393 that this is achieved with a moderate computation budget, which can be parallelised over multiple
394 machines.
395

396 **Accuracy** We compare the final structures to the experimentally determined structures to measure
397 their accuracy using metrics such as TM-score, GDT_TS¹⁷ and RMSD. All of these accuracy mea-
398 sures require geometric alignment between the candidate structure and the experimental structure.
399 An alternative accuracy measure which requires no alignment is the Local Distance Difference
400 Test (IDDT¹⁸) which measures the percentage of native pairwise distances D_{ij} under 15 Å, with
401 sequence offsets $\geq r$ residues, that are realised in a candidate structure (as d_{ij}) within a tolerance

*While there is no guarantee that these potentials have equivalent scale, scaling parameters on the terms were introduced and chosen by cross-validation on CASP12 FM domains. In practice equal weighting for all terms was found to lead to the best results.

402 of the true value, averaging across tolerances of 0.5, 1, 2 and 4 Å (without stereochemical checks).

$$lDDT_r = \frac{100}{4L} \sum_{t \in \{0.5, 1, 2, 4\}} \sum_{i=1}^L \frac{\sum_{j, |i-j| \geq r, D_{ij} < 15} \mathbb{1}(|D_{ij} - d_{ij}| < t)}{\sum_{j, |i-j| \geq r, D_{ij} < 15} 1}. \quad (5)$$

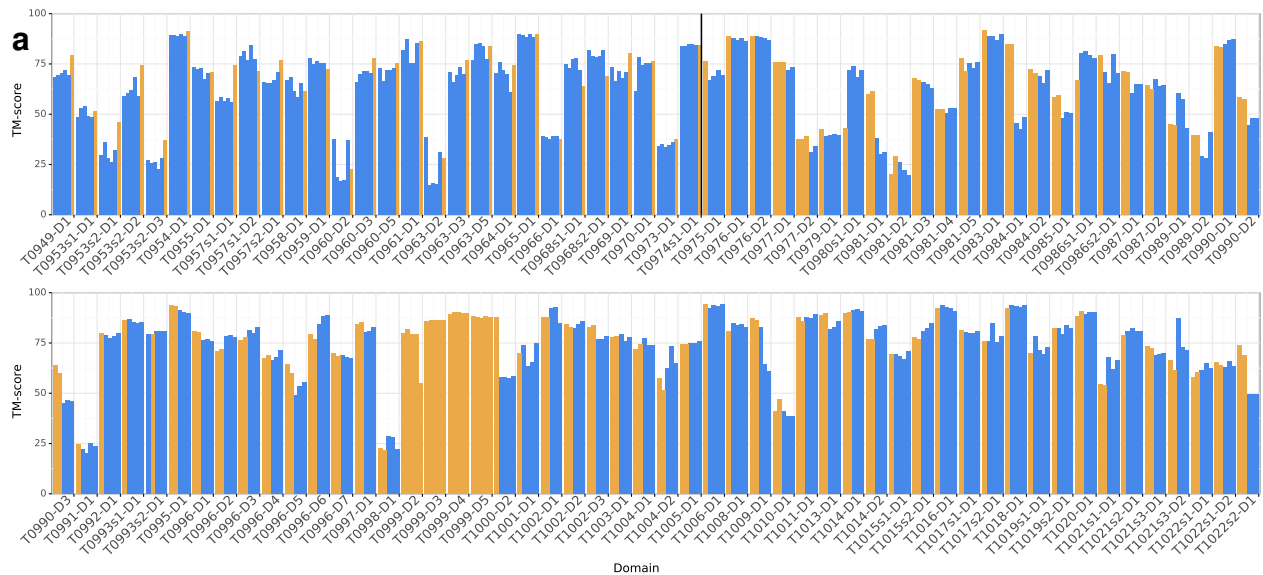
403 Since the distogram predicts pairwise distances, we can introduce *distogram lDDT* (DLDDT),
404 a measure like lDDT computed directly from the distograms’ probabilities.

$$DLDDT_r = \frac{100}{4L} \sum_{t \in \{0.5, 1, 2, 4\}} \sum_{i=1}^L \frac{\sum_{j, |i-j| \geq r, D_{ij} < 15} P(|D_{ij} - d_{ij}| < t \mid \mathcal{S}, \text{MSA}(\mathcal{S}))}{\sum_{j, |i-j| \geq r, D_{ij} < 15} 1} \quad (6)$$

405 Since distances between residues nearby in the sequence are often short, easier to predict and
406 are not critical in determining the overall fold topology, we set $r = 12$, considering only those
407 distances for residues with a sequence separation ≥ 12 . Since we predict C_β distances, for this
408 work we compute both lDDT and DLDDT using the C_β distances.

409 **Full chains without domain segmentation** Parameterising proteins of length L by two torsion
410 angles per residue, the dimension of space of structures grows as $2L$, so searching for structures
411 of large proteins becomes much harder. Traditionally this problem is addressed by splitting longer
412 protein chains into pieces, termed domains, which fold independently. However, the problem of
413 domain segmentation from the sequence alone is itself difficult and error-prone. For this work, we
414 avoided domain segmentation and folded entire chains. Typically multiple sequence alignments
415 are based upon a given domain segmentation, but we used a *sliding windows* approach, computing
416 a full-chain multiple sequence alignment to predict a baseline full-sequence distogram. We then
417 compute MSAs for subsequences of the chain, trying windows of size 64, 128, 256 with offsets
418 at multiples of 64. Each of these MSAs gives rise to an individual distogram corresponding to an
419 on-diagonal square of the full-chain distogram. We average all these distograms together, weighted
420 by the number of sequences in the MSA to produce an average full-chain distogram which is more
421 accurate in regions where many alignments can be found.

422 **CASP13 results** In CASP13 the 5 AlphaFold submissions were from 3 different systems, all us-
423 ing potentials based on the neural network distance predictions. Before T0975, two systems based
424 on simulated annealing and fragment assembly (and using 40 bin distance distributions) were used.
425 From T0975 on, newly trained 64-bin distogram predictions were used and structures were gen-
426 erated by the gradient descent system described here (3 independent runs) as well as one of the
427 fragment assembly systems (5 independent runs). 5 submissions were chosen from these 8 struc-
428 tures (the lowest potential structure generated by each independent run) with the first submission
429 (‘top-1’) being the lowest-potential structure generated by gradient descent. The remaining four
430 submissions were the four best other structures, with the fifth being a gradient descent structure if
431 none had been chosen for position 2, 3 or 4. All submissions for T0999 were generated by gradient
432 descent. Figure 10a shows the methods used for each submission, comparing with ‘back-fill’ struc-
433 tures generated by a single run of gradient descent for targets before T0975. Table 10b shows that



Method	FM	TBM/FM	TBM	All
Top-1	58.0	68.1	76.2	69.9
Best-of-5	62.6	73.6	78.6	73.2
1× gradient descent	58.4	71.6	76.3	70.4
1× fragment assembly	54.3	69.9	74.5	68.0

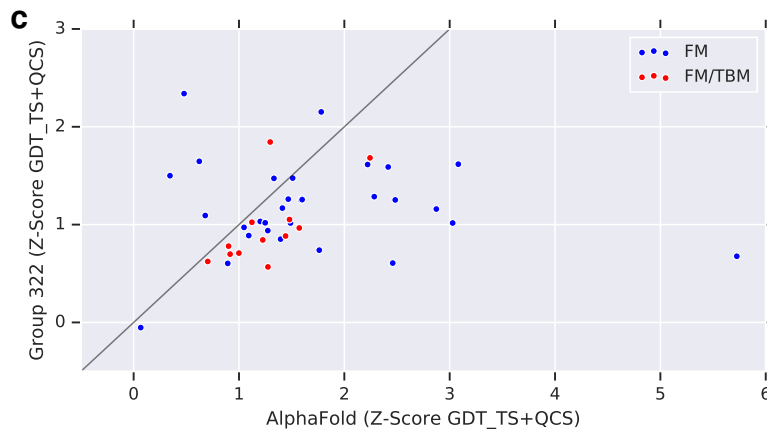


Fig. 10 | AlphaFold CASP13 results. (a) The TM-score for each of the 5 AlphaFold CASP13 submissions are shown. Simulated annealing with fragment assembly entries are shown in blue. Gradient-descent entries are shown in yellow. Gradient descent was only deployed for targets T0975 and later, so to the left of the black line we also show the results for a single, ‘back-fill’, run of gradient descent for each earlier target using the deployed system. T0999 (1589 residues) was manually segmented based on HHpred¹⁹ homology matching. (b) Average TM-scores of the AlphaFold CASP13 submissions ($n = 104$ domains), comparing the first model submitted, the best-of-5 model (submission with highest GDT), a single run of full-chain gradient descent (a CASP13 run for T0975 and later, back-fill for earlier targets) and a single CASP13 run of fragment assembly with domain segmentation (using a gradient descent submission for T0999). (c) Assessors’ formula standardised (z) scores of $GDT_{TS} + QCS$ ²⁰, best-of-5 for CASP FM ($n = 31$) and FM/TBM ($n = 12$) domains comparing AlphaFold with the closest competitor (group 322), coloured by domain category.

434 the gradient descent method deployed later in CASP performed better than the fragment assembly
435 method, in each category. Figure 10c compares the accuracy of the AlphaFold submissions for FM
436 and FM/TBM domains with the next best group 322. For the CASP13 assessment full chains were
437 relaxed with Rosetta relax with a potential of $V_{\text{Talaris2014}} + 0.2V_{\text{distance}}$ (weighting determined by
438 cross-validation) and submissions from all the systems were ranked based on this potential.

439 **Biological relevance** There is a wide range of uses of predicted structures, all with different ac-
440 curacy requirements, from generally understanding the fold shape to understanding detailed side-
441 chain configurations in binding regions. Contact predictions alone can guide biological insight²¹,
442 for instance targeting mutations to destabilise the protein. The accuracy of the contact predictions
443 shown in Table 1c indicates that the AlphaFold contact predictions exceed the state of the art. Here
444 we present further results which indicate that AlphaFold’s accuracy improvements lead to more ac-
445 curate interpretation of function; better interface prediction for protein-protein interaction; better
446 binding pocket prediction and improved molecular replacement in crystallography.

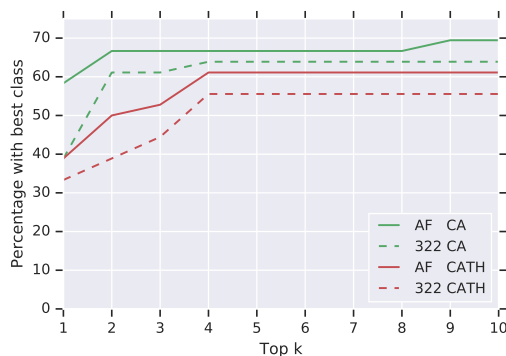


Fig. 11 | Correct fold identification by structural search in CATH. For each of the FM or TBM/FM domains, the top-1 submission and ground-truth are compared to all 30 744 CATH S40 non-redundant domains with TM-align²². For the 36 domains where there is a good ground-truth match (score > 0.5), we show the percentage of decoys where a domain with the same CATH code (in red, CA in green. CAT results are close to CATH results) as the top ground-truth match is in the at-most top-k matches with score > 0.5. Curves are shown for AlphaFold and the next-best group (322). AlphaFold predictions determine the matching fold more accurately. Determination of the matching CATH domain can give insight into the function of a new protein.

447 Often protein function can be inferred by finding homologous proteins of known function.
448 Figure 11 shows that AlphaFold’s FM predictions give greater accuracy in structure-based search
449 for homologous domains in the CATH database.

450 Protein-protein interaction is an important domain for understanding protein function that
451 has hitherto largely been limited to template-based models because of the need for high accuracy
452 predictions, though there has been moderate success²³ in docking with predicted structures up
453 to 6 Å RMSD. Figure 12 shows that AlphaFold’s predictions improve accuracy in the interface

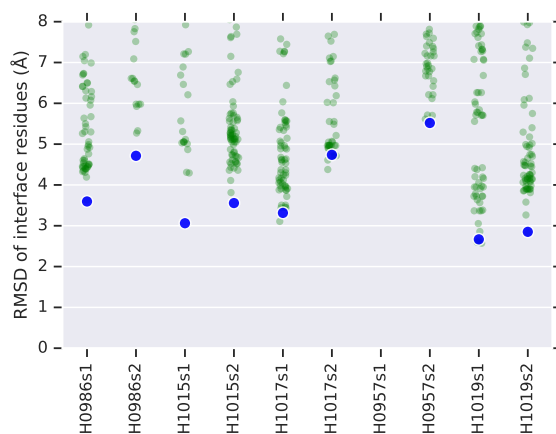


Fig. 12 | Accuracy of predictions for interfaces. For the five all-groups heterodimer CASP13 targets the full-atom RMSDs of the interface residues (residues with a ground truth inter-chain heavy atom distance $< 10 \text{ \AA}$) are computed for all groups' chain submissions, relative to the target complex. Results $> 8 \text{ \AA}$ are not shown. AlphaFold achieves consistently high accuracy interface regions and for 4 out of 5 targets predicts both chains' interfaces below $< 5 \text{ \AA}$.

454 regions of chains in hetero-dimer structures and are likely better candidates for docking, though
 455 docking did not form part of the AlphaFold system and all submissions were for isolated chains
 456 rather than complexes.

457 Further evidence of AlphaFold reaching accuracy sufficient for biological relevance is shown
 458 in Figure 13. The images show the pocket in T1011 indicating that the accuracy gain in Al-
 459 phaFold's structure prediction can lead to more accurate prediction of pocket geometry and thus
 460 the binding of ligands.

461 So far only template-based predictions have been able to deliver the most accurate predic-
 462 tions. While AlphaFold is able to match template-based modelling without using templates, and
 463 in some cases outperform other methods (e.g. T0981-D5, 72.8 GDT_TS, and T0957s1-D2, 88.0
 464 GDT_TS, two TBM-hard domains where AlphaFold's top-1 model is 12 GDT_TS better than any
 465 other top-1 submission) accuracy for FM targets still lags that for TBM targets and can still not
 466 be relied upon for detailed understanding of hard structures. In an analysis of the performance of
 467 CASP13 TBM predictions for Molecular Replacement, Read et al.²⁵ reported that the AlphaFold
 468 predictions (raw coordinates, without B-factors) led to a marginally greater log-likelihood gain
 469 (LLG) than those of any other group, indicating that these improved structures can assist in phas-
 470 ing for X-ray crystallography.

471 **Interpretation of distogram neural network** We have shown that the deep distance prediction
 472 neural network achieves high accuracy, but we would like to understand how the network arrives at
 473 its distance predictions, and in particular to understand how the inputs to the model affect the final



Fig. 13 | Ligand pocket visualizations for T1011 PDB 6M9T: EP3 receptor bound to misoprostol-FA²⁴ (a) the native structure showing the ligand in a pocket. (b) AlphaFold’s submission 5 (78.0 GDT_TS) made without knowledge of the ligand shows a pocket more similar to the true pocket than that of (c) the best other submission (322 model 3, 68.7 GDT_TS). Both submissions are aligned to the native using the same subset of residues from the helices close to the ligand pocket and visualized with the interior pocket together with the native ligand position.

474 prediction. This might lead to understanding of the folding mechanisms or suggest improvements
 475 to the model. However, deep neural networks are complex non-linear functions of their inputs, and
 476 so this attribution problem is difficult, under-specified and an on-going topic of research. Even
 477 so, there are a number of methods for such analysis: here we apply Integrated Gradients [26] to
 478 our trained distogram network to indicate the location of input features which affect the network’s
 479 predictions of a particular distance.

480 Given the expected value of the distance between any two residues I and J , $d^{I,J}(x)$, we can
 481 consider its derivatives with respect to the input features $x_{i,j,c}$, where i and j are residue indices and
 482 c is the feature channel index. The attribution function, as calculated using Integrated Gradients,
 483 of the expected distance between residues I and J with respect to the input features is then defined
 484 as

$$S_{i,j,c}^{I,J} = (x_{i,j,c} - x'_c) \int_{\alpha=0}^1 d\alpha \frac{\partial d^{I,J}(\alpha x + (1 - \alpha)x')}{\partial x_{i,j,c}}, \quad (7)$$

$$\text{s.t. } \sum_{i,j,c} S_{i,j,c}^{I,J} = d^{I,J}(x) - d^{I,J}(x'), \quad (8)$$

485 where x' is a reference set of features; in this case we average the input features spatially:

$$x'_c = \frac{1}{N^2} \sum_{i=0, j=0}^{N,N} x_{i,j,c}. \quad (9)$$

486 The derivatives of d can be calculated using backpropagation on the trained distogram network,
 487 and the integral over α is approximated as a numerical summation.

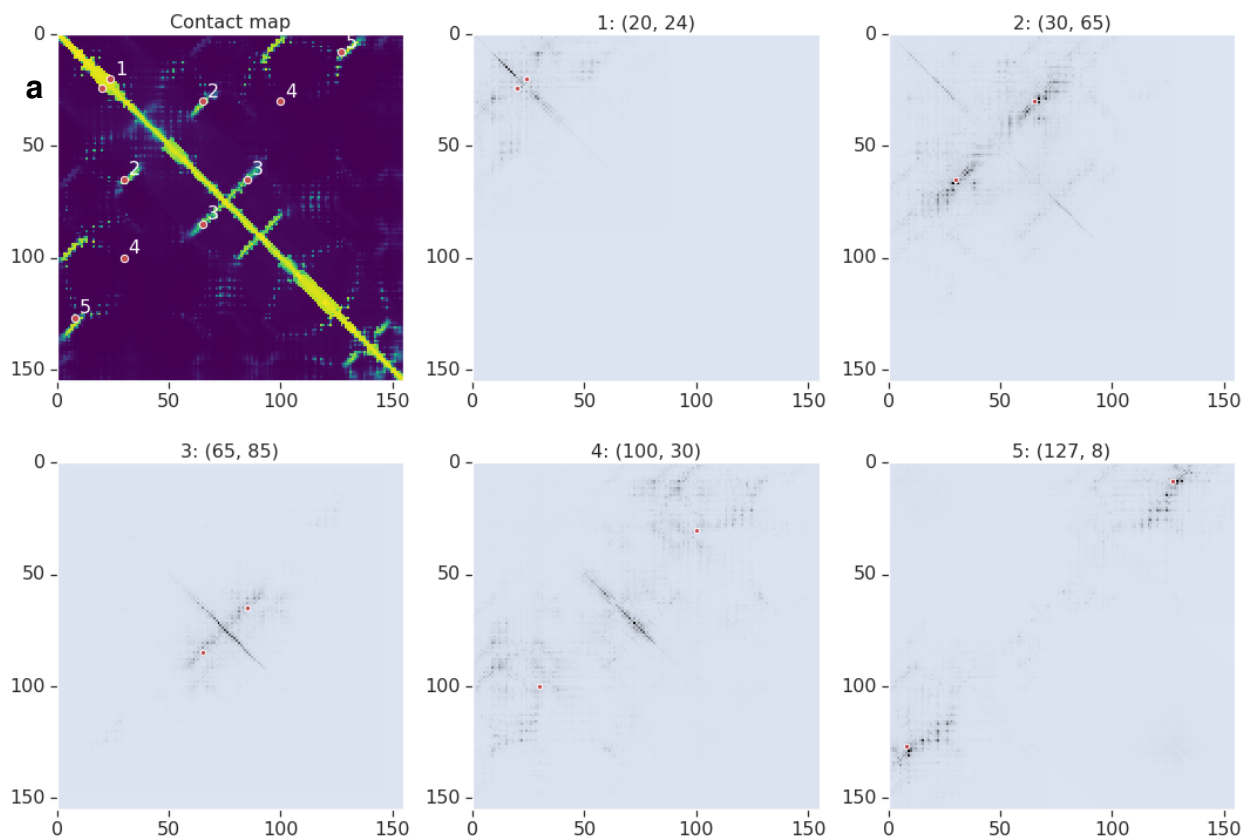


Fig. 14 | Attribution map of distogram network The contact probability map of T0986s2, and the summed absolute value of the Integrated Gradient, $\sum_c |S_{i,j,c}^{I,J}|$, of the input 2D features with respect to the expected distance between five different pairs of residues (I, J): (1) a helix self-contact, (2) a long-range stand-strand contact, (3) a medium-range strand-strand contact, (4) a non-contact and (5) a very long-range strand-strand contact. Each pair is shown as two red dots on the diagrams. Darker means higher attribution weight.

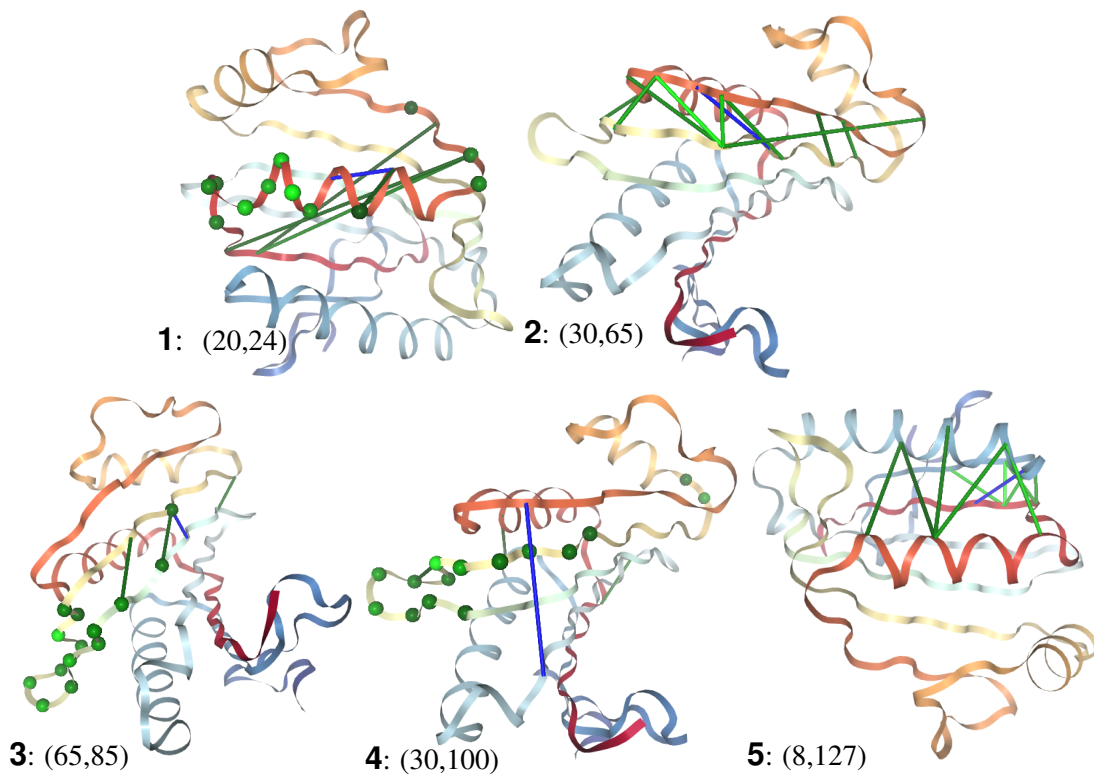


Fig. 15 | Attribution shown on predicted structure. For T0986s2 (TM-score 0.8), the top 10 input pairs, including self-pairs, with highest attribution weight for each of the five output pairs shown in Figure 14, are shown as lines (or spheres, for self-pairs) coloured by sensitivity, lighter green is more sensitive, and the output pair shown as a blue line.

488 In Figure 14, plots of summed absolute Integrated Gradient, $\sum_c |S_{i,j,c}^{I,J}|$, are shown for se-
489 lected I, J output pairs in T0986s2 and in Figure 15, the top-10 highest attribution input pairs for
490 each output pair are shown on top of AlphaFold’s top-1 predicted structure. The attribution maps
491 are sparse and highly structured, closely reflecting the predicted geometry of the protein. For the
492 four in-contact pairs presented (1, 2, 3, 5), all the highest attribution pairs are pairs within or be-
493 tween the secondary structure that one or both the output pair are members of. In (1) the helix
494 residues are important as well as connections between the strands which follow either end of the
495 helix, which might indicate strain on the helix. In (2) all the most important residue pairs connect
496 the same two strands, whereas in (3) a mix of inter-strand pairs and strand residues are most salient.
497 In (5) the most important pairs involve the packing of nearby secondary structure elements to the
498 strand and helix. For the non-contacting pair (4), the most important input pairs are the residues
499 that are geometrically between I and J in the predicted protein structure. Furthermore, most of the
500 high attribution input pairs are themselves in contact.

501 Since the network is tasked with predicting the spatial geometry, with no structure available
502 at the input, these patterns of interaction indicate that the network is using intermediate predictions
503 to discover important interactions and channelling information from related residues to refine the
504 final prediction.

505 References

- 506 1. Gilliland, G. *et al.* The Protein Data Bank. *Nucleic Acids Res.* **28**, 235–242 (2000).
- 507 2. Dawson, N. *et al.* CATH: An expanded resource to predict protein function through struc-
508 ture and sequence. *Nucleic Acids Res.* (2017).
- 509 3. Mirdita, M. *et al.* Uniclust databases of clustered and deeply annotated protein sequences and
510 alignments. *Nucleic Acids Res.* **45**, D170–D176 (2016).
- 511 4. Remmert, M., Biegert, A., Hauser, A. & Söding, J. HHblits: lightning-fast iterative protein
512 sequence searching by HMM-HMM alignment. *Nature Methods* **9**, 173 (2012).
- 513 5. Seemayer, S., Gruber, M. & Söding, J. CCMpred—fast and precise prediction of protein
514 residue–residue contacts from correlated mutations. *Bioinformatics* **30**, 3128–3130 (2014).
- 515 6. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database
516 search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).
- 517 7. Wang, S., Sun, S., Li, Z., Zhang, R. & Xu, J. Accurate De Novo Prediction of Protein Contact
518 Map by Ultra-Deep Learning Model. *PLoS Comput. Biol.* **13**, 999–1006 (2017).
- 519 8. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition. *arXiv*
520 *preprint arXiv:1512.03385* **abs/1512.03385** (2015).
- 521 9. Yu, F. & Koltun, V. Multi-scale context aggregation by dilated convolutions. *arXiv preprint*
522 *arXiv:1511.07122* (2015).

- 523 10. Oord, A. v. d. *et al.* Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*
524 (2016).
- 525 11. Clevert, D.-A., Unterthiner, T. & Hochreiter, S. Fast and Accurate Deep Network Learning
526 by Exponential Linear Units (ELUs). *arXiv preprint arXiv:1511.07289* (2015).
- 527 12. Jones, D. T. & Kandathil, S. M. High precision in protein contact prediction using fully
528 convolutional neural networks and minimal sequence features. *Bioinformatics* **1**, 8 (2018).
- 529 13. Kabsch, W. & Sander, C. Dictionary of protein secondary structure: pattern recognition of
530 hydrogen-bonded and geometrical features. *Biopolymers* **22**, 2577–2637 (1983).
- 531 14. Yang, Y. *et al.* Sixty-five years of the long march in protein secondary structure prediction:
532 the final stretch? *Briefings Bioinf.* **19**, 482–494. ISSN: 1467-5463 (2016).
- 533 15. Jones, D., Singh, T., Kosciolok, T. & Tetchner, S. MetaPSICOV: Combining coevolution
534 methods for accurate prediction of contacts and long range hydrogen bonding in proteins.
535 *Bioinformatics* **31**, 999–1006 (2015).
- 536 16. Liu, D. C. & Nocedal, J. On the limited memory BFGS method for large scale optimization.
537 *Math. Program.* **45**, 503–528 (1989).
- 538 17. Zemla, A., Venclovas, Č., Moulton, J. & Fidelis, K. Processing and analysis of CASP3 protein
539 structure predictions. *Proteins* **37**, 22–29 (1999).
- 540 18. Mariani, V., Biasini, M., Barbato, A. & Schwede, T. IDDT: a local superposition-free score
541 for comparing protein structures and models using distance difference tests. *Bioinformatics*
542 **29**, 2722–2728 (2013).
- 543 19. Söding, J., Biegert, A. & Lupas, A. The HHpred interactive server for protein homology
544 detection and structure prediction. *Nucleic acids research* **33**, W244–W248 (2005).
- 545 20. Cong, Q. *et al.* An automatic method for CASP9 free modeling structure prediction assess-
546 ment. *Bioinformatics* **27**, 3371–3378 (2011).
- 547 21. Kayikci, M. *et al.* Protein contacts atlas: visualization and analysis of non-covalent contacts
548 in biomolecules. *Nat. Struct. Mol. Biol.* **25**, 185–194 (2018).
- 549 22. Zhang, Y. & Skolnick, J. TM-align: A protein structure alignment algorithm based on TM-
550 score. *Nucleic Acids Research* **33**, 2302–2309 (2005).
- 551 23. Tovchigrechko, A., Wells, C. & Vakser, I. Docking of protein models. *Protein Sci.* **11**, 1888–
552 1896 (2002).
- 553 24. Audet, M. *et al.* Crystal structure of misoprostol bound to the labor inducer prostaglandin E
554 2 receptor. *Nature chemical biology* **15**, 11 (2019).
- 555 25. Read, R. Molecular replacement in CASP13. *In preparation* (2019).
- 556 26. Sundararajan, M., Taly, A. & Yan, Q. Axiomatic Attribution for Deep Networks. *CoRR*
557 **abs/1703.01365** (2017).

558 **2 Acknowledgements**

559 We thank Clemens Meyer for assistance in preparing the paper; Ben Coppin, Oriol Vinyals, Marek
560 Barwinski, Ruoxi Sun, Carl Elkin, Peter Dolan, Matthew Lai and Yujia Li for their contributions
561 and support; Olaf Ronneberger for reviewing the paper; the rest of the DeepMind team for their
562 support; the CASP13 organisers and the experimentalists whose structures enabled the assessment.

563 **3 Author contributions**

564 R.E., J.J., J.K., L.S., A.S., C.Q., T.G., A.Z., A.B., H.P. and K.S. designed and built the AlphaFold
565 system with advice from D.S., K.K. and D.H..

566 D.J. provided advice and guidance on protein structure prediction methodology.

567 S.P. contributed to software engineering.

568 S.C., A.N., K.K. and D.H. managed the project.

569 J.K., A.S., T.G., A.Z., A.B., R.E., P.K. and J.J. analysed the CASP results for the paper.

570 A.S., J.K. wrote the paper with contributions from J.J., R.E., L.S., T.G., A.Z., D.J., P.K., K.K. and
571 D.H.

572 A.S. led the team.

573 **4 Competing financial interests**

574 A.S., J.K., L.S., R.E., H.P., C.Q., K.S. and J.J. have filed provisional patent application 62/734,757.

575 A.S. and J.J. have filed provisional patent application 62/734,773. A.S., J.K., T.G., J.J., L.S., R.E.

576 and C.Q. have filed provisional patent application 62/770,490. J.J., A.S., R.E., A.B., T.G. and A.Z.

577 have filed provisional patent application 62/774,071. The remaining authors declare no competing

578 financial interests.

579 **5 Materials and correspondence**

580 Correspondence and materials requests should be sent to Andrew Senior: andrewsenior@google.com

581 **6 Reporting summary**

582 Further information on experimental design is available in the Nature Research Reporting Sum-
583 mary linked to this article.

584 **7 Code availability**

585 Upon publication we will make available source code for the distogram and torsion predictions,
586 together with neural network weights and input data for the CASP13 targets.

587 We make use of several open-source libraries to conduct our experiments, particularly HH-
588 blits¹, PSI-BLAST² and the machine learning framework TensorFlow[§] along with the TensorFlow
589 library Sonnet[¶] which provides implementations of individual model components³. We also used
590 Rosetta⁴ under license.

591 **8 Data availability**

592 The following public datasets were used in this work:

- 593 • PDB 2018-03-15
- 594 • CATH 2018-03-16
- 595 • Uniclust30 2017-10
- 596 • PSI-BLAST `nr` dataset (as of 2017-12-15)

597 We will make available our train/test split (CATH domain codes).

598 **9 Extended data**

599 The following tools and data set versions were used for the CASP system and for subsequent
600 experiments.

- 601 • PDB 2018-03-15
- 602 • CATH 2018-03-16
- 603 • HHblits based on version 3.0-beta.3 (3 iterations, E-value 1e-3)
- 604 • HHpred web server
- 605 • Uniclust30 2017-10
- 606 • PSI-BLAST version 2.6.0 `nr` dataset (as of 2017-12-15) (3 iterations, E-value 1e-3)
- 607 • SST web server (March 2019)
- 608 • BioPython v1.65
- 609 • Rosetta v3.5
- 610 • PyMol for structure visualisation.
- 611 • TM-align 20160521

[§]<https://github.com/tensorflow/tensorflow>

[¶]<https://github.com/deepmind/sonnet>