

# Improved protein structure prediction using predicted interresidue orientations

Jianyi Yang<sup>a,1</sup>, Ivan Anishchenko<sup>b,c,1</sup>, Hahnbeom Park<sup>b,c</sup>, Zhenling Peng<sup>d</sup>, Sergey Ovchinnikov<sup>e</sup>, and David Baker<sup>b,c,f,2</sup>

<sup>a</sup>School of Mathematical Sciences, Nankai University, 300071 Tianjin, China; <sup>b</sup>Department of Biochemistry, University of Washington, Seattle, WA 98105; <sup>c</sup>Institute for Protein Design, University of Washington, Seattle, WA 98105; <sup>d</sup>Center for Applied Mathematics, Tianjin University, 300072 Tianjin, China; <sup>e</sup>John Harvard Distinguished Science Fellowship Program, Harvard University, Cambridge, MA 02138; and <sup>f</sup>Howard Hughes Medical Institute, University of Washington, Seattle, WA 98105

Edited by William F. DeGrado, University of California, San Francisco, CA, and approved November 27, 2019 (received for review August 22, 2019)

The prediction of interresidue contacts and distances from coevolutionary data using deep learning has considerably advanced protein structure prediction. Here, we build on these advances by developing a deep residual network for predicting interresidue orientations, in addition to distances, and a Rosetta-constrained energy-minimization protocol for rapidly and accurately generating structure models guided by these restraints. In benchmark tests on 13th Community-Wide Experiment on the Critical Assessment of Techniques for Protein Structure Prediction (CASP13)and Continuous Automated Model Evaluation (CAMEO)-derived sets, the method outperforms all previously described structureprediction methods. Although trained entirely on native proteins, the network consistently assigns higher probability to de novodesigned proteins, identifying the key fold-determining residues and providing an independent quantitative measure of the "ideality" of a protein structure. The method promises to be useful for a broad range of protein structure prediction and design problems.

protein structure prediction | deep learning | protein contact prediction

lear progress in protein structure prediction was evident in the recent 13th Community-Wide Experiment on the Critical Assessment of Techniques for Protein Structure Prediction (CASP13) structure-prediction challenge (1). Multiple groups showed that application of deep learning-based methods to the protein structure-prediction problem makes it possible to generate fold-level accuracy models of proteins lacking homologs in the Protein Data Bank (PDB) (2) directly from multiple sequence alignments (MSAs) (3-6). In particular, AlphaFold (A7D) from DeepMind (7) and Xu with RaptorX (4) showed that distances between residues (not just the presence or absence of a contact) could be accurately predicted by deep learning on residue-coevolution data. The 3 top-performing groups (A7D, Zhang-Server, and RaptorX) all used deep residual-convolutional networks with dilation, with input coevolutionary coupling features derived from MSAs, either using pseudolikelihood or by covariance matrix inversion. Because these deep learningbased methods produce more complete and accurate predicted distance information, 3-dimensional (3D) structures can be generated by direct optimization. For example, Xu (4) used Crystallography and NMR System (CNS) (8) and the Alpha-Fold group (7) used gradient descent following conversion of the predicted distances into smooth restraints. Progress was also evident in protein structure refinement at CASP13 using energy-guided refinement (9-11).

In this work, we integrate and build upon the CASP13 advances. Through extension of deep learning-based prediction to interresidue orientations in addition to distances, and the development of a Rosetta-based optimization method that supplements the predicted restraints with components of the Rosetta energy function, we show that still more accurate models can be generated. We also explore applications of the model to the protein design problem. To facilitate further development in this rapidly moving field, we make all of the codes for the improved method available.

# **Results and Discussion**

**Overview of the Method.** The key components of our method (named transform-restrained Rosetta [trRosetta]) include 1) a deep residual-convolutional network which takes an MSA as the input and outputs information on the relative distances and orientations of all residue pairs in the protein and 2) a fast Rosetta model building protocol based on restrained minimization with distance and orientation restraints derived from the network outputs.

**Predicting interresidue geometries from MSAs using a deep neural network.** Unlike most other approaches to contact/distance predictions from MSAs, in addition to  $C_{\beta}$ – $C_{\beta}$  distances, we also sought to predict interresidue orientations (Fig. 1*A*). Orientations between residues 1 and 2 are represented by 3 dihedral ( $\omega$ ,  $\theta_{12}$ ,  $\theta_{21}$ ) and 2 planar angles ( $\varphi_{12}$ ,  $\varphi_{21}$ ), as shown in Fig. 1*A*. The  $\omega$  dihedral measures rotation along the virtual axis connecting the  $C_{\beta}$  atoms of the 2 residues, and  $\theta_{12}$ ,  $\varphi_{12}$  ( $\theta_{21}$ ,  $\varphi_{21}$ ) angles specify the direction of the  $C_{\beta}$  atom of residue 2 (1) in a reference frame centered on residue 1 (2). Unlike *d* and  $\omega$ ,  $\theta$  and  $\varphi$  coordinates are asymmetric and depend on the order of residues (1–2 and 2–1 pairs yield different coordinates, which is the reason why the  $\theta$  and  $\varphi$  maps in *SI Appendix*, Fig. S1 are asymmetric). Together, the 6 parameters *d*,  $\omega$ ,  $\theta_{12}$ ,  $\varphi_{12}$ ,  $\theta_{21}$ , and  $\varphi_{21}$  fully define the relative positions of the backbone atoms of 2 residues. All of the

# Significance

Protein structure prediction is a longstanding challenge in computational biology. Through extension of deep learningbased prediction to interresidue orientations in addition to distances, and the development of a constrained optimization by Rosetta, we show that more accurate models can be generated. Results on a set of 18 de novo-designed proteins suggests the proposed method should be directly applicable to current challenges in de novo protein design.

Author contributions: D.B. designed research; J.Y., I.A., H.P., Z.P., and S.O. performed research; J.Y., I.A., H.P., Z.P., S.O., and D.B. analyzed data; and J.Y., I.A., and D.B. wrote the paper.

The authors declare no competing interest.

This article is a PNAS Direct Submission.

Published under the PNAS license.

This article contains supporting information online at https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1914677117/-/DCSupplemental.

Data deposition: The multiple sequence alignments for proteins in the benchmark datasets, the codes for interresidue geometries prediction, and the Rosetta protocol for restraint-guided structure generation discussed in this paper are available at https:// yanglab.nankai.edu.cn/trRosetta/ and https://github.com/gjoni/trRosetta.

<sup>&</sup>lt;sup>1</sup>J.Y. and I.A. contributed equally to this work.

<sup>&</sup>lt;sup>2</sup>To whom correspondence may be addressed. Email: dabaker@uw.edu.



Fig. 1. Predicting interresidue geometries and protein 3D structure from a multiple sequence alignment. (A) Representation of the rigid-body transform from one residue to another using angles and distances. (B) Architecture of the deep neural network with multiobjective training to predict interresidue geometries from an MSA. (C) Outline of the structure-modeling protocol based on the restraints derived from the predicted distance and orientation (see *Methods* for details).

coordinates show characteristic patterns (*SI Appendix*, Fig. S1), and we hypothesized that a deep neural network could be trained to predict these.

The overall architecture of the network is similar to those recently described for distance and contact prediction (3, 4, 7, 12). Following RaptorX-Contact (4, 12) and AlphaFold (7), we learn probability distributions over distances and extend this to orientation features. The central part of the network is a stack of dilated residual-convolutional blocks that gradually transforms 1- and 2-site features derived from the MSA of the target to predict interresidue geometries for residue pairs (Fig. 1B) with  $C_{\beta}$  atoms closer than 20 Å. The distance range (2 to 20 Å) is binned into 36 equally spaced segments, 0.5 Å each, plus one bin indicating that residues are not in contact. After the last convolutional layer, the softmax function is applied to estimate the probability for each of these bins. Similarly,  $\omega$ ,  $\theta$  dihedrals and  $\varphi$ angle are binned into 24, 24, and 12, respectively, with 15° segments (+ one no-contact bin) and are predicted by separate branches of the network. Branching takes place at the very top of the network, with each branch consisting of a single convolutional layer followed by softmax. The premise for such hard parameter sharing at the downstream layers of the networks is that correlations between the different objectives (i.e., orientations and distance) may be learned by the network, potentially yielding better predictions for the individual features. We used cross-entropy to measure the loss for all branches; the total loss is the sum over the 4 per-branch losses with equal weight. Previous work (4) implicitly captured some orientation information by predicting multiple interresidue distances ( $C_{\beta}$ - $C_{\beta}$ ,  $C_{\alpha}$ - $C_{\alpha}$ ,  $C_{\alpha}$ - $C_{g}$ ,  $C_g-C_g$ , and N-O), but in contrast to our multitask-learning approach, a separate network was used for each of the objectives. Our network was trained on a nonredundant (at 30% sequence identity) dataset from PDB consisting of 15,051 proteins (structure release dates before 1 May 2018). The trained network is available for download at https://github.com/gjoni/trRosetta.

We couple the derivation of residue–residue couplings from MSAs by covariance matrix inversion to the network by making the former part of the computation graph in TensorFlow (13).

Sequence reweighting, calculation of one-site amino acid frequencies, entropies, and coevolutionary couplings and related scores take place on the GPU, and the extracted features are passed into the convolutional layers of the network (most previous approaches have precomputed these terms). We took advantage of our recent observation (14) that with proper regularization, covariance matrix inversion yields interresidue couplings (Methods) with only minor decrease in accuracy compared to pseudolikelihood approaches like GREMLIN (15) (the latter are prohibitively slow for direct integration into the network). Since the MSA-processing steps are now cheap to compute (compared to the forward and backward passes through the network during parameter training), this coupled network architecture allows for data augmentation by MSA subsampling during training. At each training epoch, we use a randomly selected subset of sequences from each original MSA, so that each time the network operates on different inputs.

Structure modeling from predicted interresidue geometries. Following AlphaFold, we generated 3D structures from the predicted distances and orientations using constrained minimization (Fig. 1C). Discrete probability distributions over the predicted orientation and distance bins were converted into interresidue interaction potentials by normalizing all of the probabilities by the corresponding probability at the last bin (Methods) and smoothing using the spline function in Rosetta. These distance- and orientationdependent potentials were used as restraints, together with the Rosetta centroid level (coarse-grained) energy function (16), and folded structures satisfying the restraints were generated starting from conformations with randomly selected backbone dihedral angles by 3 rounds of quasi-Newton minimization within Rosetta. Only short-range (sequence separation <12) restraints were included in the first round; medium-range (sequence separation <24) restraints were added in the second round, and all were included in the third. A total of 150 coarse-grained models were generated using different sets of restraints obtained by selecting different probability thresholds for inclusion of the predicted distances and orientations in modeling.

The 50 lowest-energy backbone + centroid models were then subjected to Rosetta full-atom relaxation, including the distance and orientation restraints, to add in side chains and make the structures physically plausible. The lowest-energy full-atom model was then selected as the final model. The structure generation protocol is implemented in PyRosetta (17) and is available as a web server at https://yanglab.nankai.edu.cn/trRosetta/.

# Benchmark Tests on CASP13 and Continuous Automated Model Evaluation Datasets.

Accuracy of predicted interresidue geometries. We tested the performance of our network on 31 free-modeling (FM) targets from CASP13. (None of these were included in the training set, which is based on a pre-CASP PDB set.) The precision of the derived contacts, defined as the fraction of top L/n (n = 1, 2, 5) predicted contacts realized in the native structure, is summarized in Table 1 and SI Appendix, Table S1. For the highest probability 7.5% of the distance/orientation predictions (Fig. 2C), there is a good correlation between modes of the predicted distance/orientation distributions and the observed values (Fig. 2C): Pearson r for distances is 0.72, and circular correlation  $r_c$  (18) for  $\omega$ ,  $\theta$ , and  $\varphi$ are 0.62, 0.77, and 0.60, respectively. The predicted probability of the top L long- + medium-range contacts correlates well (r = 0.84) with their actual precision (Fig. 2B). This correlation between predicted probability and actual precision allows us to further improve the results by feeding a variety of MSAs generated with different e-value cutoffs or originating from searches against different databases, into the network and selecting the one that generates predictions with the highest predicted accuracy.

**Comparison with baseline network.** We evaluated our extensions to previous approaches by generating a baseline model to predict distances only, with no MSA subsampling and selection; the contact prediction accuracy of this network is comparable to previously described models (3, 12, 19, 20). Incorporating MSA subsampling during training and extending the network to also predict interresidue orientations improve contact prediction accuracy by 1.7 and 2.2%, respectively. Subsequent alignment selection improves performance an additional 3.1% on the CASP13 FM set (Table 1, last row). The improvements described above, together with increasing the number of layers in the network, increase the accuracy of predicted contacts by 7.6% over the baseline network on the CASP13 FM set. Although we ensured that

Table 1. Precision (%) of the top *L* predicted contacts on CASP13 and CAMEO targets

Method	CASP13 FM domains		CAMEO very hard targets	
	<i>s</i> ≥ 24	<i>s</i> ≥ 12	<i>s</i> ≥ 24	<i>s</i> ≥ 12
RaptorX-Contact	44.7	61.3	NA	NA
TripleRes	42.3	60.9	NA	NA
trRosetta	51.9	70.2	48.0	62.8
Baseline*	44.3	60.7	41.6	57.5
Baseline+1 <sup>+</sup>	46.0	62.2	43.1	57.4
Baseline+1+2 <sup>‡</sup>	48.2	64.6	44.4	58.7
Baseline+1+2+3 <sup>§</sup>	51.3	69.3	46.1	61.4

The values for other methods are slightly different from those listed on the CASP13 website (http://predictioncenter.org/casp13/), probably due to different treatment of target length *L* (i.e., length of full sequence or length of domain structures; the latter is used here). The sequence separation between 2 residues *i* and *j* is denoted by  $s (=|i\cdotj|)$ .

\*Baseline trRosetta model consists of 36 residual blocks and was trained without MSA subsampling or selection to predict distances only.

<sup>†</sup>1: adding MSA subsampling during training.

<sup>‡</sup>2: extending the network to predict orientations.

 $^{\S}3$ : MSA selection based on predicted probability of the top L long- + medium-range contacts.

there is no overlap between the training and test sets by selecting pre-CASP PDBs only (before 1 May 2018), our model was trained at a later date when more sequences were available; we also included metagenomic sequence data. Hence, we may be overestimating the gap in performance between our method and those used by other groups in CASP13; future blind tests in CASP will be important in confirming these improvements. Nevertheless, the gain in performance with respect to the baseline model is independent of the possible variations in the training sets and sequence databases. All of the targets in the Continuous Automated Model Evaluation (CAMEO) validation set below are more recent than both structural and sequence data in the training set. *Accuracy of predicted structure models.* We tested our method on the CASP13 FM targets, with results shown in Fig. 3. The average TM-score (21) of our method is 0.625, which is 27.3% higher

TM-score (21) of our method is 0.625, which is 27.3% higher than that (0.491) by the top Server group Zhang-Server (Fig. 3A). Our method also outperforms the top Human group A7D by 6.5% (0.625 vs. 0.587; Fig. 3B). The relatively poor performance on T1021s3-D2 (the outlier in the upper triangle of Fig. 3B) reflects the MSA-generation procedure: the majority of sequence homologs in the full-length MSA for T1021S3 only covers the first of the 2 domains; performance is significantly improved (TM-score increased from 0.38 to 0.63; the TM-score of A7D model is 0.672) using a domain-specific MSA. An example of the improved performance of our method is shown in Fig. 3C for the CASP13 target T0950; the TM-score of this model is 0.716, while the highest values obtained during CASP13 are: RaptorX-DeepModeller (0.56), BAKER-ROSETTASERVER (0.46), Zhang-Server (0.44), and A7D (0.43).

Fig. 3A deconstructs the contributions to the improved performance of the different components of our approach. When modeling is only guided by the distance predictions from the baseline network (no orientations and no MSA subsampling and selection; "baseline" bar in Fig. 2A), the TM-score is 0.537, lower than A7D but significantly higher than Zhang-Server and RaptorX. When predicted distances from the complete network are used, the TM-score increases to 0.592, higher than that of A7D. When the orientation distributions are included, the TM-score is further increased to 0.625. The folding is driven by restraints; very similar models are generated without the Rosetta centroid terms, and very poor models are generated without the restraints. To compare our Rosetta minimization protocol (trRosetta) to CNS (8), we obtained predicted distance restraints and structure models for all CASP13 FM targets from the RaptorX-Contacts server, which uses CNS for structure modeling (4), and used the distance restraints to generate models with trRosetta. The average TM-score of the trRosetta models is 0.45 compared to 0.36 for the RaptorX CNS models; the improvement is likely due to both improved sampling and the supplementation of the distance information with the general Rosetta centroid energy function.

Comparison between distance and orientation-based folding. Both predicted distance and orientation can guide folding alone. The average TM-score of coarse-grained models for the CASP13 FM targets is 0.57 when folding with predicted orientation alone and 0.55 when folding with predicted distance only (SI Appendix, Fig. S24). Relaxation improved the TM-score to 0.58 and 0.59 for orientation and distance guided folding, respectively (SI Ap*pendix*, Fig. S2B). The differences in quality of models generated using either source of information alone suggest that the 2 are complementary, and indeed better models are generated using both distance and orientation information (SI Appendix, Fig. S2). Validation on hard targets from the CAMEO experiments. We further tested our method on 131 hard targets from the CAMEO experiments (22) over the 6 mo between 8 December 2018 and 1 June 2019. The results for contact prediction are summarized in Table 1 and Fig. 24; as in the case of the CASP13 targets, our method improves over the baseline network. The results for



**Fig. 2.** Accuracy of predicted interresidue geometries. (*A*) Contribution of different factors to the increase in trRosetta performance on CASP13's free modeling and CAMEO's very hard targets. Incorporation of MSA subsampling, orientations, and MSA selection in the modeling pipeline increases precision of the top *L* long-range predicted contacts by 1.7% (red bar), 2.2% (yellow), and 3.1% (green), respectively, and increasing the depth of the network from 36 to 61 residual blocks boosts the performance by an additional 0.6% (orange bar). (*B*) Correlation between predicted probability of the top *L* long- + medium-range contacts and their actual precision measured based on the native structures. (C) Distribution of predicted probabilities for residue pairs to be within 20 Å in the native structure; populations in blue and red correspond to residue pairs with  $d \le 20$  Å and d > 20 Å in experimental structures, respectively. (D) Correlations between actual rigid-body transform parameters from the experimental structures with the modes of the predicted distributions for the most reliable long- and medium-range contacts from the top 7.5% percentile; color coding indicates probability density.

structure modeling are shown in Fig. 3D. The contributions of different components to our method are presented in *SI Appendix*, Fig. S4. On these targets, the average TM-score of our method is 0.621, which is 8.9 and 24.7% higher than Robetta and HHpredB, respectively. We note that the definition of "hard" is looser than the CASP definition; a hard target from CAMEO can have close templates in PDB. Making the definition of "hard" more stringent by requiring the TM-score of the HHpredB server to be less than 0.5 reduces the number of targets to 66. On this harder set, the TM-score for our method is 0.534, 22% higher than the top server Robetta and 63.8% higher than the baseline server HHpredB. Fig. 3E shows an example of a CAMEO target where our method predicts very accurate models (5WB4\_H). For this target, the TMscores of the template-based models by HHpredB, IntFOLD5-TS, and RaptorX are about 0.4. In comparison, the TM-score of our predicted model is 0.921, which is also higher than the top server Robetta (0.879).

Accuracy estimation for predicted structure models. We sought to predict the TM-score of the final structure model using the 131 hard targets from CAMEO. We found that, unlike direct couplingbased methods such as GREMLIN, the depth of the MSA did not have a good correlation with the accuracy of the derived contacts. Instead, a high correlation (Pearson r = 0.90) between the average probability of the top-predicted contacts and the actual precision was observed (SI Appendix, Fig. S3A). The average contact probability also correlates well with the TM-score of the final structure models (r = 0.71; *SI Appendix*, Fig. S3B). To obtain a structure-based accuracy metric, we rerelaxed the top 10 models without any restraints. The average pairwise TM-score between these 10 nonconstrained models also correlates with the TM-score of the final models (r = 0.65; *SI Appendix*, Fig. S3C). Linear regression against the average contact probability and the extent of structural displacement without the restraints gave a quite good correlation between predicted and actual TM-score (r = 0.84; SI Appendix, Fig. S3D). We used this method to provide an estimated model accuracy.

**Refinement of predicted models.** As noted above, CASP13 showed that protein structure-refinement methods can consistently improve models for cases where the sampling problem is more tractable

(smaller monomeric proteins). We first evaluated the iterative hybridization protocol (23) previously used to improve models generated using direct contacts predicted from GREMLIN on the entire set of CASP13 and CAMEO targets (SI Appendix, Fig. S5). Incorporating our network-derived distance predictions resulted in consistent improvement in model quality when the starting model's TM-score was over 0.7, in a few cases by more than 10%in TM-score. We also tested the incorporation of the networkderived distance restraints into the more compute-intensive structure refinement protocol we used in CASP13 (10) on the CASP13 FM targets with an estimated starting TM-score >0.6 that were not heavily intertwined oligomers and not bigger than 250 residues. Consistent improvements were observed on a set of 6 such targets (SI Appendix, Fig. S6), with an average TM-score improvement of about 4%. The net improvement in prediction for these targets using the combination of our structuregeneration method and refinement using the distance predictions is indicated by the red points in Fig. 3B.

Assessing the Ideality of de Novo Protein Designs. Following up on the AlphaFold group's excellent CASP13 prediction of the designed protein T1008, we systematically compared the ability of trRosetta to predict the structure of de novo-designed proteins from single sequences compared to native proteins in the same length range. We collected a set of 18 de novo-designed proteins of various topologies (24–26) ( $\alpha$ ,  $\beta$ , and  $\alpha/\beta$ ) with coordinates in the PDB and a set of 79 natural proteins of similar size selected from the CAMEO set and ran the trRosetta protocol to predict interresidue geometries (Fig. 4A) and 3D models (Fig. 4B; examples of 3D models are in Fig. 4 C-E). There is a clear difference in performance for natural proteins and de novo designs: the latter are considerably more accurate. The predicted structures of the designed proteins are nearly superimposable on the crystal structures, which is remarkable given that there is no coevolution information whatsoever for these computed sequences, which are unrelated to any naturally occurring protein.

The high-accuracy structure prediction in the absence of coevolutionary signal suggests the model is capturing fundamental features of protein sequence–structure relationships. To further



**Fig. 3.** Comparison of model accuracy. (*A*) Average TM-score of all methods on the 31 FM targets of CASP13. The colored stacked bar indicates the contributions of different components to our method. A7D was the top human group in CASP 13; Zhang-Server and RaptorX were the top 2 server groups. (*B*) Head-to-head comparison between our method and the A7D's TM-scores over the 31 FM targets (blue points; red points are for 6 targets with extensive refinement). (*C*) Structures for the CASP13 target T0950; the native structure and the predicted model are shown in gray and rainbow cartoons, respectively. (*D*) Comparison between our method and the top servers from the CAMEO experiments. (*E*) Native structure (in gray) and the predicted model (in rainbow) for CAMEO target 5WB4\_H. In all of these comparisons, it should be emphasized that the CASP and CAMEO predictions, unlike ours, were made blindly.

investigate this, we performed an exhaustive mutational scanning of the "wild-type" sequences for 3 designs of distinct topology (24–26) (Fig. 4 *C*–*E* and *SI Appendix*, Fig. S7). For each single amino acid substitution at each position, we calculated the change in the probability of the top *L* long- + medium-range contacts  $[-\log(P_{mutant}/P_{WT})]$ . Mutations of core hydrophobic residues and of glycine residues in the  $\beta$ -turns produced large decreases in the probability of the designed structure. The effects of mutations depend strongly on context: the substitutions of the same amino acid type at different positions produce quite different changes in probability (Fig. 4 *C*–*E*), which go far beyond the averaged out information provided by simple BLOSUM and PAM.

# Discussion

The results presented here suggest that the orientation information predicted from coevolution can improve structure prediction. Tests on the CASP13 and CAMEO sets suggest that our combined method outperforms all previously described methods, as it should, as we have attempted to build on the many advances made by many groups in CASP13. However, it should be emphasized that retrospective analyses such as those carried out in this paper are no substitute for blind prediction experiments (as in the actual CASP13 and CAMEO) and that future CASP and CAMEO testing will be essential. Although not fully explored in this work, the integrated network architecture allows for backpropagation of gradients down to the MSA-processing step, making it possible to learn optimal sequence reweighting and regularization parameters directly from data rather than using manually tuned values. To enable facile exploration of the ideas presented in this paper and in CASP13, the codes for the orientation prediction from coevolution data and the Rosetta protocol for structure generation from predicted distances and orientations are all available at https://yanglab.nankai.edu.cn/trRosetta/ and https://github.com/gjoni/trRosetta.

The accurate prediction of the structure of de novo-designed proteins in the complete absence of coevolutionary signal has implications for both the model and protein design generally. First, the model is clearly learning general features of protein structures. This is not surprising given that the direct couplings derived by the coevolutionary analysis on a protein family are the 2-body terms in a generative model for the sequences in the family, and thus training on these couplings for a large number of protein families is equivalent to training on large sets of protein sequences for each structure in the training set. From the design point of view, we have asserted previously that de novo-designed proteins are "ideal" versions of naturally occurring proteins (27); the higher probability assigned by the model to designed proteins compared to naturally occurring proteins makes this assertion quantitative. Remarkably, similar "ideal" features appear to have been distilled from native protein analysis by expert protein designers to be incorporated into designed proteins, and extracted by deep learning in the absence of any expert intervention. Our finding that the model provides information on the contribution of each amino acid in a designed protein to the determination of the fold by the sequence suggests the model should be directly applicable to current challenges in de novo protein design.

SIOPHYSICS AND COMPUTATIONAL BIOLOG



**Fig. 4.** trRosetta accurately predicts structures of de novo-designed proteins and captures effects of mutations. Differences in the accuracy of predicted contacts (*A*) and trRosetta models (*B*) for de novo-designed (blue) and naturally occurring (orange) proteins of similar size from single amino acid sequences. (*C–E*) Examples of trRosetta models for de novo designs of various topology:  $\beta$ -barrel, PDB ID 6DOT (*C*);  $\alpha$ -helical IL2-mimetic, PDB ID 6DC6 (*D*); and Foldit design with  $\alpha/\beta$  topology, PDB ID 6MRS (*E*). Experimental structures are in gray, and models are in rainbow. Frames show experimental structures color-coded by estimated tolerance to single-site mutations (red, less tolerant; blue, more tolerant); the 8 residues least tolerant to mutation are in stick representation, and glycine residues are indicated by arrows. Heat maps on the right show the change in probability of the designed fold for substitutions of the same residue type (indicated at top) at different sequence positions (indicated at bottom).

This work also demonstrates the power of modern deep learning packages such as TensorFlow in making deep learning model development accessible to nonexperts. The distance and orientation prediction method described here performs comparably or better than models previously developed by leading experts (of course we had the benefit of their experience), despite the relative lack of expertise with deep learning in our laboratory. These packages have now opened up deep learning to scientists generally-the challenge is more to identify appropriate problems, datasets and features than to formulate and train the models. The method developed here is immediately applicable to problems ranging from cryoEM model fitting to sequence generation and structure optimization for de novo protein design.

#### Methods

#### Benchmark Datasets.

Training set for the neural network. To train the neural network for the prediction of distance and orientation distributions, a training set consisting of 15,051 protein chains was collected from the PDB. First, we collected 94,962 X-ray entries with resolution  $\leq 2.5$  Å (PDB snapshot as of May first 2018), then extracted all protein chains with at least 40 residues, and finally removed redundancy at 30% sequence identity cutoff, resulting in a set of 16,047 protein chains with the average length of 250 amino acids. All of the corresponding primary sequences were then used as queries to collect MSAs using the iterative procedure described below. Only chains with at least 100 sequence homologs in the MSA were selected for the final training set.

*Independent test sets.* Two independent test sets are used to test our method. The first is the 31 FM domains (25 targets) from CASP13 (first target released on 1 May 2018). The second one is from the CAMEO experiment. We collected 131 CAMEO hard targets released between 8 December 2018 and 1 June 2019, along with all of the models submitted by public servers during this period. Note that for the

MSA generation and selection. The precision of predicted distance and orientation distribution usually depends on the availability of an MSA with 'good' quality. A deep MSA is usually preferable but not always better than a shallow MSA (see the examples provided in ref. 3). In this work, 5 alternative alignments are generated for each target. The first 4 are generated independently by searching the Uniclust30 database (version 2018\_08) with HHblits (version 3.0.3) (28) with default parameters at 4 different e-value cutoffs:  $1e^{-40}$ ,  $1e^{-10}$ ,  $1e^{-3}$ , and 1. The last alignment was generated by several rounds of iterative HHblits searches with gradually relaxed e-value cutoffs  $(1e^{-80}, 1e^{-70}, \dots, 1e^{-10}, 1e^{-8}, 1e^{-6}, and 1e^{-4})$ , followed by the hmmsearch (version 3.1b2) (29) against the metagenome sequence database (20) in case not enough sequences were collected at previous steps. The metagenome database includes about 7 billion protein sequences from the following resources: 1) JGI Metagenomes (7,835 sets), Metatranscriptomes (2,623 sets), and Eukaryotes (891 genomes); 2) UniRef100; 3) NCBI TSA (2,616 sets); and 4) genomes manually collected from various genomic centers and online depositories (2,815 genomes). To avoid attracting distant homologs at early stages and making alignment unnecessarily deep, the search was stopped whenever either of the 2 criteria were met: at least 2,000 sequences with 75% coverage or 5,000 sequences with 50% coverage (both at 90% sequence identity cutoff) were collected. The final MSAs for the test datasets are available at https://yanglab.nankai.edu.cn/trRosetta/.

CASP13 dataset, the full protein sequences rather than the domain sequences are

used in all stages of our method to mimic the situation of the CASP experiments.

### Interresidue Geometries Prediction by Deep Residual Neural Networks.

**Protein structure representation.** In addition to the traditional interresidue distance matrices, we also make use of orientation information to make the representation locally informative. For a residue pair (*i*, *j*), we introduce  $\omega$  dihedral between  $C_{\alpha}$ ,  $C_{\beta}$  of one residue and  $C_{\beta}$ ,  $C_{\alpha}$  of the other, as well as 2 sets of spherical coordinates centered at each of the residues and pointing to the  $C_{\beta}$  atom of the other residue. These 6 coordinates (*d*,  $\omega$ ,  $\theta_{ji}$ ,  $\varphi_{ji}$ ,  $\theta_{ji}$ ,  $\varphi_{ji}$ ) are sufficient to fully define the relative orientation of 2 residues with respect to one another. Additionally,

as described below, any biasing energy term defined along these coordinates can be straightforwardly incorporated as restraints in Rosetta.

*Input features.* All of the input features for the network are derived directly from the MSA and are calculated on-the-fly. The 1D features include: 1) one-hot-encoded amino acid sequence of the query protein (20 feature channels), 2) position-specific frequency matrix (21 features: 20 amino acids + 1 gap), and 3) positional entropy (1 feature). These 1D features are tiled horizontally and vertically and then stacked together to yield  $2 \times 42 = 84$  2D feature maps.

Additionally, we extract pair statistics from the MSA. It is represented by couplings derived from the inverse of the shrunk covariance matrix constructed from the input MSA. First we compute 1-site and 2-site frequency counts  $f_i(A) = \frac{1}{M_{eff}} \sum_{m=1}^{M} w_m \delta_{A,A_{i,m}}$  and  $f_{i,j}(A, B) = \frac{1}{M_{eff}} \sum_{m=1}^{M} w_m \delta_{A,A_{i,m}}$ , where A and B denote amino acid identities (20 + gap),  $\delta$  is the Kronecker delta, indices i, j run through columns in the alignment, and the summation is over all M sequence in the MSA;  $w_m$  is the inverse of the number of sequences in the MSA, which share at least 80% sequence identity with sequence m (including itself);  $M_{eff} = \sum_{m=1}^{M} w_m$ . We then calculate the sample covariance matrix

$$C_{i,j}^{A,B} = f_{i,j}(A,B) - f_i(A)f_j(B)$$
 [1]

and find its inverse (also called the precision matrix) after shrinkage (i.e., regularization by putting additional constant weights on the diagonal):

$$S_{i,j}^{A,B} = \left(C_{i,j}^{A,B} + \frac{4.5}{\sqrt{M_{eff}}}\delta_{i,j}\delta_{A,B}\right)^{-1}.$$
 [2]

(More details on tuning the regularization weight in Eq. **2** are provided in *SI* Appendix, Fig. S8). The 21 × 21 coupling matrices  $s_{i,j}^{A,B}$  of the precision matrix (Eq. **2**) are flattened, and the resulting  $L \times L \times 441$  feature matrix contributes to the input of the network. The above couplings (Eq. **2**) are also converted into single values by computing their Frobenius norm for nongap entries:

$$\mathbf{s}_{i,j}^{*} = \sqrt{\sum_{A=1}^{20} \sum_{B=1}^{20} \left( \mathbf{s}_{i,j}^{A,B} \right)^{2}},$$
 [3]

followed by the average product correction (APC):

$$s_{i,j} = s_{i,j}^* - s_{i,j}^* s_{i,.}^* / s_{i,.}^*$$
 [4]

where  $s_{i,j}^*$ ,  $s_{i,i,j}^*$ , and  $s_{i,j}^*$  are row, column, and full averages of the  $s_{i,j}^*$  matrix, respectively. The coefficient in Eq. **2** was manually tuned on a nonredundant set of 1,000 proteins to maximize accuracy of the top *L* predicted contacts. From our experience, the final results are quite stable to the particular choice of the regularization coefficient in Eq. **2**. To summarize, the input tensor has 526 feature channels: 84 (transformed 1D features) + 441 (couplings; Eq. **2**) + 1 (APC score; Eq. **4**).

Network architecture. The network takes the above  $L\times L\times 526$  tensor as the input and applies a sequence of 2D convolutions to simultaneously predict 4 objectives: 1 distance histogram (*d* coordinate) and 3 angle histograms ( $\omega$ ,  $\theta$  and  $\varphi$  coordinates). After the first layer, which transforms the number of input features down to 64 (2D convolution with filter size 1), the stack of 61 basic residual blocks with dilations are applied. Dilations cycle through 1, 2, 4, 8, and 16 (12 full cycles in total). After the last residual block, the network branches out into 4 independent paths—one per objective—with each path consisting of a single 2D convolution followed by softmax activation. Since maps for *d* and  $\omega$  coordinates are symmetric, we enforce symmetry in the network right before the corresponding 2 branches by adding transposed and untransposed feature maps from the previous layer. All convolution operations, except the first and the last, use 64 3 × 3 filters; ELU activations are applied throughout the network.

Training. We use categorical cross-entropy to measure the loss for all 4 objectives. The total loss is the sum over the 4 individual losses with equal weight (= 1.0), assuming that all coordinates are equally important for structure modeling. During training, we randomly subsample the input MSAs, uniformly in the log scale of the alignment size. Big proteins of more than 300 amino acids long are randomly sliced to fit 300 residue limit. Each training epoch runs through the whole training set, and 100 epochs are performed in total. Adam optimizer with the learning rate  $1e^{-4}$  is used. All trainable parameters are restrained by the  $L_2$  penalty with the  $1e^{-4}$  weight. Dropout keeping probability 85% is used. We train 5 networks with random 95/5% training/validation splits and use the average over the 5 networks as the final prediction. Training a single network takes ~9 d on one NVIDIA Titan RTX GPU.

# Structure Determination by Energy Minimization with Predicted Restraints.

*Converting distance and orientation distribution to energy potential.* The major steps for structure modeling from predicted distributions are shown in Fig. 1C. For each pair of residues, the predicted distributions are converted into energy potential following the idea of Dfire (30). For the distance distribution, the probability value for the last bin, i.e., [19.5, 20], is used as a reference state to convert the probability values into scores by the following equation:

$$\operatorname{core}^{d}(i) = -\ln(p_{i}) + \ln\left(\left(\frac{d_{i}}{d_{N}}\right)^{a} p_{N}\right), i = 1, 2, \cdots, N,$$
[5]

where  $p_i$  is the probability for the *i*th distance bin, *N* is the total number of bins,  $\alpha$  is a constant (= 1.57) for distance-based normalization, and  $d_i$  is the distance for the *i*th distance bin. For the orientation distributions, the conversion is similar but without normalization, i.e.,

s

S

core<sup>o</sup>(i) = 
$$-\ln(p_i) + \ln(p_N), i = 1, 2, \dots, N.$$
 [6]

All scores are then converted into smooth energy potential by the *spline* function in Rosetta and used as restraints to guide the energy minimization. The range for distances is [0, 20 Å] with a bin size of 0.5 Å, while for orientations, the ranges are  $[0, 360^\circ]$  for  $\theta$  and  $\omega$ , and  $[0, 180^\circ]$  for  $\varphi$ , all with a bin size of 15°; corresponding cubic spline curves are generated from the discrete scores defined by Eqs. 5 and 6. For the distance-based potential, the *AtomPair* restraint is applied. For the  $\theta$ - and  $\omega$ -based potential, the *Dihedral* restraint is applied.

*Quasi-Newton-based energy minimization and full atom-based refinement.* To speed up the modeling, coarse-grained (centroid) models are first built with the quasi-Newton-based energy minimization (*MinMover*) in Rosetta. A centroid model is a reduced representation of protein structure, in which the backbone remains fully atomic but each side chain is represented by a single artificial atom (centroid). The optimization is based on the L-BFGS algorithm (*lbfgs\_armijo\_nonmonotone*). A maximum of 1,000 iterations is used, and the convergence cutoff is 0.0001. Besides the restraints introduced above, the following Rosetta energy terms are also used: ramachandran (rama), the omega and the steric repulsion van der Waals forces (vdw), and the centroid backbone hydrogen bonding (cen\_hb). More details about these energy terms can be found in ref. 16. The weights for the *AtomPair*, *Dihedral*, and *Angle* restraints, rama, omega, vdw, and cen\_hb, are 5, 4, 4, 1, 0.5, and 1, respectively. The final models are selected based on the total score which includes both Rosetta energy and restraints scores.

The *MinMover* algorithm is deterministic but can be easily trapped into local minima. It is sensitive to the initial structure and restraints. Two strategies are proposed to introduce randomization effect, and those models trapped into local minima can be discarded based on total energy. The first strategy is to use different starting structures with random backbone torsion angles (10 are tried). The second strategy consists of using different sets of restraints. For each residue pair, we only select a subset of restraints with probability higher than a specified threshold (from 0.05 to 0.5, with a step of 0.1).

For each starting structure, 3 different models are built by selecting different subsets of restraints based on sequence separation *s*: short range ( $1 \le s < 12$ ), medium range ( $12 \le s < 24$ ), and long range ( $s \ge 24$ ). The first one is progressively built with short-, medium-, and long-range restraints. The second one is built with short- + medium-range restraints and then with long-range restraints. The last one is built by using all restraints together.

In total, 150 (=  $10 \times 5 \times 3$ ) centroid models were generated. The top 10 models (ranked by total energy) at each of the probability cutoff are selected for full-atom relax by *FastRelax* in Rosetta. In this relax, the restraints at probability threshold 0.15 are used together with the ref2015 scoring function. The weights for the *AtomPair*, *Dihedral*, and *Angle* restraints are 4, 1, and 1, respectively.

Data Availability. The multiple sequence alignments for proteins in the benchmark datasets, the codes for interresidue geometries prediction, and the Rosetta protocol for restraint-guided structure generation are available at https://yanglab.nankai.edu.cn/trRosetta/ and https://github.com/gjoni/trRosetta.

ACKNOWLEDGMENTS. We thank Frank DiMaio and David Kim for helpful discussions. This work was supported by National Natural Science Foundation of China Grants NSFC 11871290 (to J.Y.) and 61873185 (to Z.P.); Fok Ying-Tong Education Foundation Grant 161003 (to J.Y.); Key Laboratory for Medical Data Analysis and Statistical Research of Tianjin (J.Y.); the Thousand Youth Talents Plan of China (J.Y.); the China Scholarship Council (J.Y. and Z.P.); Fundamental Research Funds for the Central Universities (to J.Y.); National Institute of General Medical Sciences Grant R01-GM092802-07 (to D.B.); National Institute of Allergy and Infectious Diseases Contract HHSN272201700059C (to D.B.); the Schmidt Family Foundation (D.B.); and

- L. A. Abriata, G. E. Tamò, M. Dal Peraro, A further leap of improvement in tertiary structure prediction in CASP13 prompts new routes for future assessments. *Proteins* 87, 1100–1112 (2019).
- 2. H. M. Berman et al., The protein data bank. Nucleic Acids Res. 28, 235-242 (2000).
- S. M. Kandathil, J. G. Greener, D. T. Jones, Prediction of interresidue contacts with DeepMetaPSICOV in CASP13. Proteins 87, 1092–1099 (2019).
- J. Xu, Distance-based protein folding powered by deep learning. Proc. Natl. Acad. Sci. U.S.A. 116, 16856–16865 (2019).
- J. Hou, T. Wu, R. Cao, J. Cheng, Protein tertiary structure modeling driven by deep learning and contact distance prediction in CASP13. Proteins 87, 1165–1178 (2019).
- W. Zheng et al., Deep-learning contact-map guided protein structure prediction in CASP13. Proteins 87, 1149–1164 (2019).
- J. R. Evans et al., "De novo structure prediction with deep-learning based scoring" in Thirteenth Critical Assessment of Techniques for Protein Structure Prediction (Protein Structure Prediction Center, 2018), pp. 1–4.
- A. T. Brünger et al., Crystallography & NMR system: A new software suite for macromolecular structure determination. Acta Crystallogr. D Biol. Crystallogr. 54, 905–921 (1998).
- L. Heo, C. F. Arbour, M. Feig, Driven to near-experimental accuracy by refinement via molecular dynamics simulations. *Proteins* 87, 1263–1275 (2019).
- H. Park et al., High-accuracy refinement using Rosetta in CASP13. Proteins 87, 1276– 1282 (2019)
- R. J. Read, M. D. Sammito, A. Kryshtafovych, T. I. Croll, Evaluation of model refinement in CASP13. Proteins 87, 1249–1262 (2019).
- S. Wang, S. Sun, Z. Li, R. Zhang, J. Xu, Accurate de novo prediction of protein contact map by ultra-deep learning model. *PLoS Comput. Biol.* 13, e1005324 (2017).
- M. Abadi et al., Tensorflow: Large-scale machine learning on heterogeneous distributed systems. arXiv:1603.04467 (14 March 2016).
- 14. J. Dauparas et al., Unified framework for modeling multivariate distributions in biological sequences. arXiv:1906.02598 (6 June 2019).
- H. Kamisetty, S. Ovchinnikov, D. Baker, Assessing the utility of coevolution-based residue-residue contact predictions in a sequence- and structure-rich era. Proc. Natl. Acad. Sci. U.S.A. 110, 15674–15679 (2013).

Office of the Director of the National Institutes of Health Grant DP5OD026389 (to S.O.).

- C. A. Rohl, C. E. Strauss, K. M. Misura, D. Baker, Protein structure prediction using Rosetta. *Methods Enzymol.* 383, 66–93 (2004).
- S. Chaudhury, S. Lyskov, J. J. Gray, PyRosetta: A script-based interface for implementing molecular modeling algorithms using Rosetta. *Bioinformatics* 26, 689–691 (2010).
- S. R. Jammalamadaka, A. Sengupta, *Topics in Circular Statistics* (World Scientific, 2001).
- Y. Li, J. Hu, C. Zhang, D. J. Yu, Y. Zhang, ResPRE: High-accuracy protein contact prediction by coupling precision matrix with deep residual neural networks. *Bio*informatics 35, 4647–4655 (2019).
- Q. Wu et al., Protein contact prediction using metagenome sequence data and residual neural networks. Bioinformatics btz477 (2019).
- Y. Zhang, J. Skolnick, Scoring function for automated assessment of protein structure template quality. *Proteins* 57, 702–710 (2004).
- J. Haas et al., Continuous Automated Model EvaluatiOn (CAMEO) complementing the critical assessment of structure prediction in CASP12. Proteins 86 (suppl. 1), 387–398 (2018).
- S. Ovchinnikov et al., Protein structure determination using metagenome sequence data. Science 355, 294–298 (2017).
- J. Dou et al., De novo design of a fluorescence-activating β-barrel. Nature 561, 485– 491 (2018).
- B. Koepnick et al., De novo protein design by citizen scientists. Nature 570, 390–394 (2019).
- D. A. Silva et al., De novo design of potent and selective mimics of IL-2 and IL-15. Nature 565, 186–191 (2019).
- N. Koga et al., Principles for designing ideal protein structures. Nature 491, 222–227 (2012).
- M. Remmert, A. Biegert, A. Hauser, J. Söding, HHblits: Lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nat. Methods* 9, 173–175 (2011).
- S. C. Potter et al., HMMER web server: 2018 update. Nucleic Acids Res. 46, W200– W204 (2018).
- H. Zhou, Y. Zhou, Distance-scaled, finite ideal-gas reference state improves structurederived potentials of mean force for structure selection and stability prediction. *Protein Sci.* 11, 2714–2726 (2002).