

# 1 **Template-based prediction of protein structure with deep learning**

2

3 Haicang Zhang<sup>1,#</sup>, Yufeng Shen<sup>1, 2, 3, 4 #</sup>

4

- 5 1. Department of Systems Biology, Columbia University, New York, NY, USA
- 6 2. Department of Biomedical Informatics, Columbia University, New York, NY, USA
- 7 3. JP Sulzberger Columbia Genome Center, Columbia University, New York, NY, USA
- 8 4. Program in Mathematical Genomics, Columbia University, New York, NY, USA

9

10 # Correspondence should be addressed to H.Z. ([hz2529@cumc.columbia.edu](mailto:hz2529@cumc.columbia.edu)) and Y.S.

11 ([ys2411@cumc.columbia.edu](mailto:ys2411@cumc.columbia.edu))

12

13

## 14 **Abstract**

15 Accurate prediction of protein structure is fundamentally important to understand biological  
16 function of proteins. Template-based modeling, including protein threading and homology  
17 modeling, is a popular method for protein tertiary structure prediction. However, accurate  
18 template-query alignment and template selection are still very challenging, especially for the  
19 proteins with only distant homologs available. We propose a new template-based modelling  
20 method called ThreaderAI to improve protein tertiary structure prediction. ThreaderAI  
21 formulates the task of aligning query sequence with template as the classical pixel classification  
22 problem in computer vision and naturally applies deep residual neural network in prediction.  
23 ThreaderAI first employs deep learning to predict residue-residue aligning probability matrix by  
24 integrating sequence profile, predicted sequential structural features, and predicted residue-  
25 residue contacts, and then builds template-query alignment by applying a dynamic programming  
26 algorithm on the probability matrix. We evaluated our methods both in generating accurate  
27 template-query alignment and protein threading. Experimental results show that ThreaderAI  
28 outperforms currently popular template-based modelling methods HHpred, CNFpred, and the  
29 latest contact-assisted method CEthreader, especially on the proteins that do not have close  
30 homologs with known structures. In particular, in terms of alignment accuracy measured with  
31 TM-score, ThreaderAI outperforms HHpred, CNFpred, and CEthreader by 56%, 13%, and 11%,  
32 respectively, on template-query pairs at the similarity of fold level from SCOPe data. And on  
33 CASP13's TBM-hard data, ThreaderAI outperforms HHpred, CNFpred, and CEthreader by 16%,

34 9% and 8% in terms of TM-score, respectively. These results demonstrate that with the help of  
35 deep learning, ThreaderAI can significantly improve the accuracy of template-based structure  
36 prediction, especially for distant-homology proteins.

37

38 **Availability:** <https://github.com/ShenLab/ThreaderAI>

39

40 **Keywords:** protein structure prediction · protein threading · deep learning · deep residual neural  
41 network

42

### 43 **1 Introduction**

44 Protein structure is fundamentally important to understand protein functions. Computational  
45 protein structure prediction remains one of the most challenging problems in structural  
46 bioinformatics. Recent progress in protein structure prediction showed that with the help of deep  
47 learning, it's possible for free modelling (FM) methods to generate fold-level accuracy models of  
48 proteins lacking homologs in protein structure library<sup>1-4</sup>. Meanwhile, as both  
49 protein sequence and structure databases expand, template-based modelling (TBM) methods  
50 remain to be very popular and useful<sup>5-7</sup> for the proteins with homologs available in protein  
51 structure library. TBM method predicts the structure of query protein by modifying the structural  
52 framework of its homologous protein with known structure in accordance with template-query  
53 alignment. The quality of TBM prediction critically relies on template-query alignment and  
54 template selection. It remains to be very challenging for TBM methods to predict structures  
55 accurately when only remote homologs which are conserved in structure but share low sequence  
56 similarity with query are available in structure library<sup>5-7</sup>.

57

58 The model accuracy of TBM method critically depends on protein features and the scoring  
59 functions that integrate these features. For protein features, sequence profiles, and protein  
60 secondary structures are widely used by exiting popular TBM methods such as HHpred<sup>8</sup>,  
61 CNFpred<sup>9</sup>, and Sparks-X<sup>10</sup>. As a result of recent progress in residue-residue contact prediction,  
62 contact information has been integrated by several recently developed methods such as  
63 DeepThreader<sup>5</sup>, CEthreader<sup>6</sup>, and EigenThreader<sup>11</sup>. For scoring functions, HHpred, Sparks-X,  
64 CEthreader, and several other methods used linear functions, while non-linear models such as

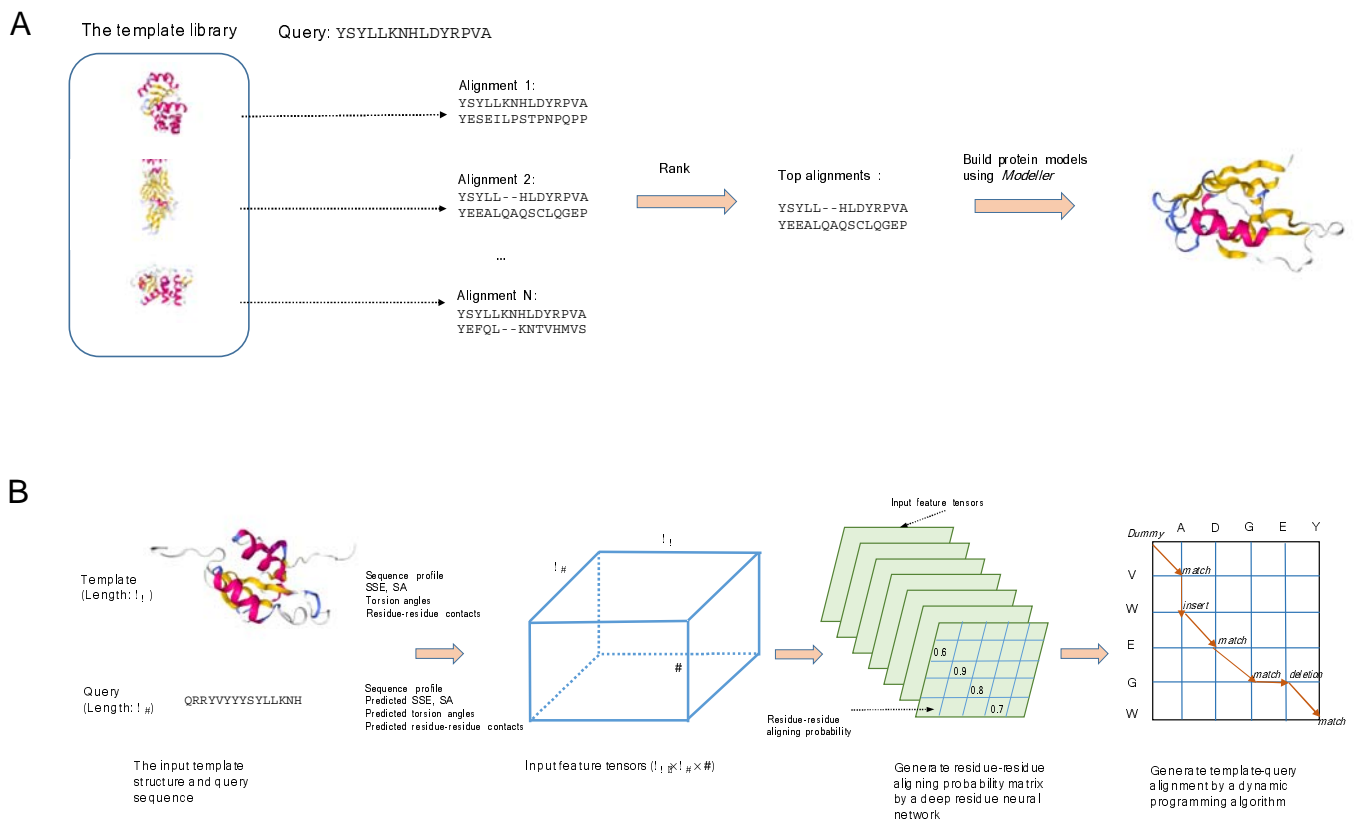
65 Random Forest model in Boost-Threader<sup>12</sup> and one-layer dense neural network in CNFPred have  
 66 shown their advantages over linear models. Inspired by the success of non-linear models in TBM  
 67 methods, we would like to study if we can improve TBM methods' model accuracy using more  
 68 advanced neural network architecture such as deep residual network which has proven very  
 69 successful in protein residue-residue contacts prediction.

70

71 In this paper, we present a new method, called ThreaderAI, which uses a deep residual neural  
 72 network to perform template-query alignment. More specifically, we formulate template-query  
 73 alignment problem as the classical pixel classification problem in computation vision. We first  
 74 adapt the deep residual neural network model to predict residue-residue aligning scoring matrix,  
 75 and then we employ a dynamic programming algorithm on the predicted scoring matrix to  
 76 generate the optimal template-query alignment.

77

## 78 2 Materials and Methods



**Figure 1.** Overview of ThreaderAI. A. The procedure of protein structure prediction using ThreaderAI. B. The procedure of aligning query with template using a deep residual neural network model and a dynamic programming algorithm.

## 79 **2.1 Overview of the method**

80 For a query protein, ThreaderAI predicts its tertiary structure through the following steps (Figure  
81 1A). First, query protein is aligned to each template in the structure library using a deep residual  
82 neural network model and a dynamic programming algorithm. Second, all the alignments are  
83 ranked based on alignment scores. Third, the final tertiary structures of query are built using  
84 Modeller<sup>13</sup> based on the top-ranked alignments.

85

86 For TBM methods, the quality of query-template alignments critically determines the quality of  
87 predicted structures<sup>5,9</sup>. ThreaderAI uses a deep residue neural network model to generate  
88 template-sequence alignment (Figure 1B). First, protein features are extracted from both  
89 template and query. Second, a deep residue neural network model is used to generate residue-  
90 residue aligning probability matrix. Third, a dynamic programming algorithm is applied on the  
91 scoring matrix to generate the final template-query alignment.

92

## 93 **2.2 Protein features**

94 We included the following features as inputs for our deep residual neural network model (also  
95 see Table 1).

96 *Sequence profile* (40 features): HHblits<sup>14</sup> was used to generate the sequence profile for both  
97 template and query. The feature vector for each residue-residue aligning pair is from the  
98 concatenation of the sequence profiles of template and query.

99

100 *Sequential Structural features* (29 features): For template, we generated its 8-class secondary  
101 structure types, real-valued solvent accessibility, and backbone dihedral angles using DSSP<sup>15</sup>.  
102 We also calculated the contact numbers of template with  $C_{\alpha}$ - $C_{\alpha}$  and  $C_{\beta}$ - $C_{\beta}$  distances of 8Å as  
103 threshold. And for Glycine, we only used its  $C_{\alpha}$  coordinates. For query, we predicted its 3-class  
104 secondary structure types, real-valued solvent accessibility, backbone dihedral angles, disordered  
105 regions, and residual level interfaces using NetSurfP2<sup>16</sup>. We also predicted these sequential  
106 structural features for template. The features of a residue-residue aligning pair are the  
107 concatenation of the structural properties of these two residues.

108

109 *Residue-residue contacts* (8 features): The residue-residue contacts of template are defined as the  
 110 residue pairs with  $C_\beta-C_\beta$  distance less than  $8\text{\AA}$ . For query, we predict its contact map using  
 111 ResPRE<sup>17</sup>. The eigenvectors and eigenvalues of the residue-residue contact matrix can capture  
 112 the intrinsic properties of protein's tertiary structure and have been used as features by recently  
 113 developed threading methods<sup>6,11</sup>. Given contact matrix  $M$  of the protein, the  $i$ th residue can be  
 114 represented as  $(\sqrt{\lambda_1}v_{1i}, \sqrt{\lambda_2}v_{2i}, \dots, \sqrt{\lambda_K}v_{Ki})$  where  $\lambda_j$  and  $v_j$  is the  $j$ th eigenvalue and  
 115 eigenvector of matrix  $M$ , respectively. Here we set  $K$  as 8. Given template and query's contact  
 116 matrices  $M_T$  and  $M_Q$ , the features of the  $i$ th residue of template aligning the  $j$ th residue of query  
 117 are defined as  $(\sqrt{\lambda_1^T\lambda_1^Q}|v_{1i}^T v_{1j}^Q|, \sqrt{\lambda_2^T\lambda_2^Q}|v_{2i}^T v_{2j}^Q|, \dots, \sqrt{\lambda_K^T\lambda_K^Q}|v_{Ki}^T v_{Kj}^Q|)$ . Heuristically, we set the  
 118 sign of each eigenvector as positive. Previous methods<sup>6,11</sup> enumerated a total of  $2^K$  possible  
 119 alignments to decide the sign of each involved eigenvector which is very time-consuming and  
 120 infeasible for neural network-based models.

sequence profiles (20×2 features)	Amino acid type distribution in multiple sequence alignment (20 features for template and 20 features for query)
sequential structural features only for template (13 features)	8-class secondary structure types (8 features)
	solvent accessibility (1 feature)
	backbone dihedral angles (2 features)
	contact numbers (2 features)
predicted sequential structural features for both template and query (8×2 features)	predicted 3-class secondary structure types (3 features)
	predicted solvent accessibility (1 feature)
	predicted backbone dihedral angles (2 features)
	predicted residual level interfaces (1 feature)
residual-residual contacts (8 features)	predicted disordered regions (1 feature)
	the dot products of the corresponding elements of top 8 eigenvectors of contact matrices of template and query

**Table 1.** Protein features used in ThreaderAI.

121

122

## 123 **2.3 Neural network architecture**

124 We employed a deep residual neural network<sup>18</sup> (ResNet) model to predict residue-residue  
125 aligning probability matrix. ResNet has proven very successful in computer vision and also in  
126 structural bioinformatics. First, convolutional layers in ResNet are capable of extracting  
127 hierarchical features or spatial patterns from images or image-like data automatically. Second,  
128 the residual component in ResNet can efficiently mitigate the issue of vanishing/exploding  
129 gradients and makes it possible to train an ultra-deep neural network model on a large scale of  
130 training data.

131

132 Specifically, for a template-query pair, the input feature tensor for our neural network model has  
133 dimensions of  $L_T \times L_Q \times d$  where  $L_T$  and  $L_Q$  denotes the lengths of template and query,  
134 respectively, and  $d$  is the number of features for each residue-residue pair. And the output for our  
135 model has dimensions of  $L_T \times L_Q$  each element of which representing residue-residue aligning  
136 probability. Our model includes 16 residue blocks<sup>18</sup> each of which includes 2 convolutional  
137 layers. Each convolutional layers used 16 filters and a kernel size of  $3 \times 3$ . We used ELU<sup>19</sup> as  
138 nonlinear activation function. Sigmoid function was used as the final layer to output residue-  
139 residue aligning probabilities.

140

## 141 **2.4 Alignment labels and training loss function**

142 We built training template-query pairs from proteins with known structures. For each template-  
143 query pair in training data, we used DeepAlign<sup>20</sup> to generate its structural alignments as ground  
144 truth. For a template with a length of  $L_T$  and a query with a length of  $L_Q$ , there are  $L_T \times L_Q$   
145 residue-residue pairs in total, in which the aligned pairs in the structural alignment are labeled as  
146 positives while the others as negatives.

147

148 Instead of using binary labels directly, we weighted<sup>9</sup> the conservation of aligned residue pairs  
149 using local TM-score<sup>21</sup>. Given a structure alignment of two proteins and the corresponding  
150 superimposition, the local TM-score of an aligned residue pair  $T_i$  and  $Q_j$  is defined as follows:

$$w_{ij} = \frac{1}{1 + (d_{ij}/d_0)^2}$$

151 where  $d_{ij}$  is the distance deviation between the two aligned residues and  $d_0$  is a normalization  
152 constant depending only on protein length. The TM-score ranges from 0 to 1, with higher values  
153 indicating more highly conserved aligned positions. And for a gap in the alignment, the local  
154 TM-score  $w_{ij}$  is equal to 0.

155  
156 The labels from the structure alignments are highly imbalanced in which the ratio of negatives  
157 over positives is proportional to the lengths of template and query. To mitigate this imbalanced  
158 labeling issue, we weighted the aligned pairs in the reference alignments with the average length  
159 of template and query.

160

161 We used cross-entropy loss as our training loss function which is defined as follows:

$$\frac{1}{N} \sum_{n=1}^N \frac{1}{L_T^{(n)} L_Q^{(n)}} \sum_i^{L_T^{(n)}} \sum_j^{L_Q^{(n)}} \left[ -L^{(n)} w_{ij}^{(n)} \log p_{ij}^{(n)} - (1 - w_{ij}^{(n)}) \log (1 - p_{ij}^{(n)}) \right]$$

162 where  $N$  is the number of protein pairs in training data and  $n$  iterates over all training samples,  
163 and  $L^{(n)}$  equals  $(L_T^{(n)} + L_Q^{(n)})/2$  meaning the average length of template and query, and  $w_{ij}^{(n)}$  and  
164  $p_{ij}^{(n)}$  are residue-residue aligning probability from our neural network model and local TM-score,  
165 respectively.

166

## 167 **2.5 Training algorithm**

168 We used AdamW algorithm<sup>22</sup> to minimize the objective function with a weight decay rate of 1e-  
169 4. For the warmup stage, we increased the learning rate from 0 to 0.01 over the first 2 epochs.

170 We also decayed the learning rate to 1e-4 with a polynomial decay policy in the following 16  
171 epochs<sup>22</sup>. Early-stopping with validation error as a metric was performed during training. The

172 model architecture and training algorithm was implemented by TensorFlow2<sup>23</sup> and run on 3  
173 NVIDIA GeForce-1080 GPUs in parallel. We set training batch size as 2 and we didn't try a  
174 larger batch size due to the limited GPU memory.

175

## 176 **2.6 Maximum accuracy algorithm**

177 Given the residue-residue aligning probability matrix of  $L_T \times L_Q$  from our neural network model,  
178 we used a dynamic programming algorithm called Maximum Accuracy algorithm (MAC)<sup>8,14</sup> to  
179 generate the final template-query alignment. MAC creates the local alignment through  
180 maximizing the sum of probabilities for each residue pair to be aligned minus a penalty  $\alpha$  which  
181 can control the alignment greediness. To find the best MAC alignment path, an optimal sub-  
182 alignment score matrix  $S$  is calculated recursively using the probability  $p_{ij}$  as substitution scores:  
183

$$S_{i,j} = \max \begin{cases} p_{ij} - \alpha \\ S_{i,j-1} - \alpha/2 \\ S_{i-1,j} - \alpha/2 \\ 0 \end{cases}$$

184  
185 Then standard traceback procedure of dynamic programming<sup>24</sup> was then applied on the score  
186 matrix  $S$  to generate the optimal local alignment. We rank the template-query alignments based  
187 on the optimal alignment scores from MAC.

188

## 189 **2.7 Dealing with proteins of variable lengths**

190 Our model has an architecture of fully convolutional neural network<sup>25</sup> in which no fully-  
191 connected layers were used. As a result, the number of parameters of our model is independent  
192 of the lengths of both template and query. Hence, our model can deal with proteins of variable  
193 lengths. In particular, zero paddings were applied so that each training sample in the same  
194 minibatch has the same size. We also filtered out the padded positions when we aggregated the  
195 final training loss.

196

## 197 **2.8 Training and test data**

198 We built the training set, validation set, and independent testing set from proteins in SCOPe40.  
199 We also included CASP13 data for testing.

200

### 201 **2.8.1 Training data**

202 We prepared template-query pairs from SCOPe40<sup>26</sup>. First, for testing purpose, we excluded the  
203 domains which share larger than 25% sequence identity with the domains in CASP13 data<sup>7</sup>. Here  
204 we used MMseqs2<sup>27</sup> to evaluate sequence identity with the default E-value of 1e-3. Second, we



205 excluded families with single domains. Third, for each class of  $\alpha$ ,  $\beta$ ,  $\alpha/\beta$ , and  $\alpha + \beta$  of  
206 SCOPe40, we randomly selected 5 folds as independent testing data and the left folds as training  
207 data. The testing and training template-query pairs were generated from testing and training folds  
208 respectively.

209  
210 Template-query pairs at the similarity of fold, superfamily and family levels were generated  
211 separately. When generating family level pairs, at most 10 pairs were randomly selected for each  
212 family. And when generating superfamily and fold level pairs, for each family pairs from the  
213 same fold, we randomly selected 1 domain from each family as its representative to form pairs.  
214 And all protein pairs with TM-score less than 0.3 were excluded. Finally, we have 53734 training  
215 pairs and 2000 validation pairs from the training folds, and 3106 pairs from testing folds.

216

## 217 **2.8.2 Test data**

218 We used two test datasets to test ThreaderAI in terms of alignment accuracy and protein  
219 threading performance, respectively. For testing alignment accuracy, we used 3106 template-  
220 query pairs (denoted as SCOPe3K data) created together with training pairs and validation pairs  
221 (see section 2.8.1). The testing template-query pairs belong to different folds with training and  
222 validation pairs. The second test set consists of 61 officially-defined CASP13<sup>7</sup> target domains  
223 under the category of Template-Based Modelling (TBM). The CPSP13 TBM data are divided  
224 into two groups by difficulty level: TBM-easy (40 targets) and TBM-hard (21 targets). We used  
225

226 To test the threading performance of ThreaderAI using CASP13 TBM data , we built our  
227 template database from PDB90 in which any two proteins share less than 90% sequence identity.  
228 We only included the structures deposited before CASP13. We also excluded the structures with  
229 more than 800 amino acids and the structures with more than 50% unobserved residues. Finally,  
230 our template library includes 50099 proteins.

231

## 232 **2.9 Evaluating metrics**

### 233 **2.9.1 Evaluating alignment accuracy**

234 For a query protein and one of a candidate template from the template library, we evaluated the  
235 alignment accuracy by evaluating the quality of the structure built from this alignment. In

236 particular, for each template-query pair, we first used ThreaderAI to generate an alignment, then  
237 built a 3D structure for the query using MODELLER<sup>13</sup> based on the alignment, and finally  
238 evaluate the similarity between the predicted structure and the ground truth structure. Here, we  
239 evaluated the quality of a 3D model by GDT<sup>28</sup> and TM-score, two widely used metric for  
240 measuring the similarity of two protein structures. GDT score is calculated based on the largest  
241 set of residue-residue pairs falling in a defined distance cutoff when superposing these two  
242 structures. GDT ranges from 0 to 100, but we normalize it by 100 so that it has a scale between 0  
243 and 1. TM-score is designed to be length-independent by introducing a length-dependent  
244 normalization factor. TM-score ranges from 0 to 1 with 1 indicating the perfect model quality.

245

### 246 **2.9.2 Evaluating threading performance**

247 We evaluated threading performance by measuring the quality of 3D models built from the top-  
248 ranked templates. Specifically, for a query protein, we used ThreaderAI to generate alignments  
249 for all the templates in template library, ranked these alignments by alignment scores and then  
250 built 3D models using MODELLER from the top five alignments. Finally, we evaluated the  
251 quality of the first-ranked and the best of top five 3D models by TM-score and GDT.

252

### 253 **2.10 Compare with previously published methods**

254 We compared ThreaderAI with several widely used threading methods including HHpred<sup>8</sup>,  
255 CNFpred<sup>9</sup>, and CEthreader<sup>6</sup>, a new threading method built upon contacts predicted by ResPre<sup>17</sup>.  
256 Here, HHpred was run with the option mact 0.1, real secondary structures for template, and  
257 predicted secondary structures for query proteins. And CEthreader was run with the mode of  
258 EigenProfileAlign in which sequence profile, secondary structures, and contact maps are used.  
259 For protein threading, we used CEthreader's suggested strategy to speedup. That is, we first run  
260 CEthreader's greedy algorithm and then selected top the 1000 templates for refinement using its  
261 enumerative algorithm. DeepThreader<sup>5</sup> is another recently developed threading software in  
262 which a linear function was used to combine local potentials from CNFpred and pairwise  
263 potentials from predicted residue-residue contacts. DeepThreader's performance wasn't shown  
264 here because its package is unavailable to the public. To be fair, for all methods we used the  
265 same template database (see section 2.8.2) and used HHblits<sup>14</sup> to build sequence profiles against

266 sequence database uniclust30\_2017\_10 built before CASP13. We used HHblits' utility script to  
267 convert HHblits' profile format to BLAST's<sup>29</sup> profile format used by CNFpred.

268

## 269 **3 Results**

	Fold level		Superfamily level		Family level		All	
	TM-score	GDT	TM-score	GDT	TM-score	GDT	TM-score	GDT
ThreaderAI	0.419	0.348	0.483	0.409	0.705	0.633	0.510	0.437
HHpred	0.268	0.226	0.411	0.353	0.675	0.607	0.416	0.362
CNFpred	0.371	0.307	0.443	0.375	0.681	0.612	0.470	0.404
CEthreader	0.377	0.313	0.425	0.356	0.641	0.568	0.456	0.389

**Table 2.** Alignment accuracy measured by TM-score and GDT on SCOPe3K data

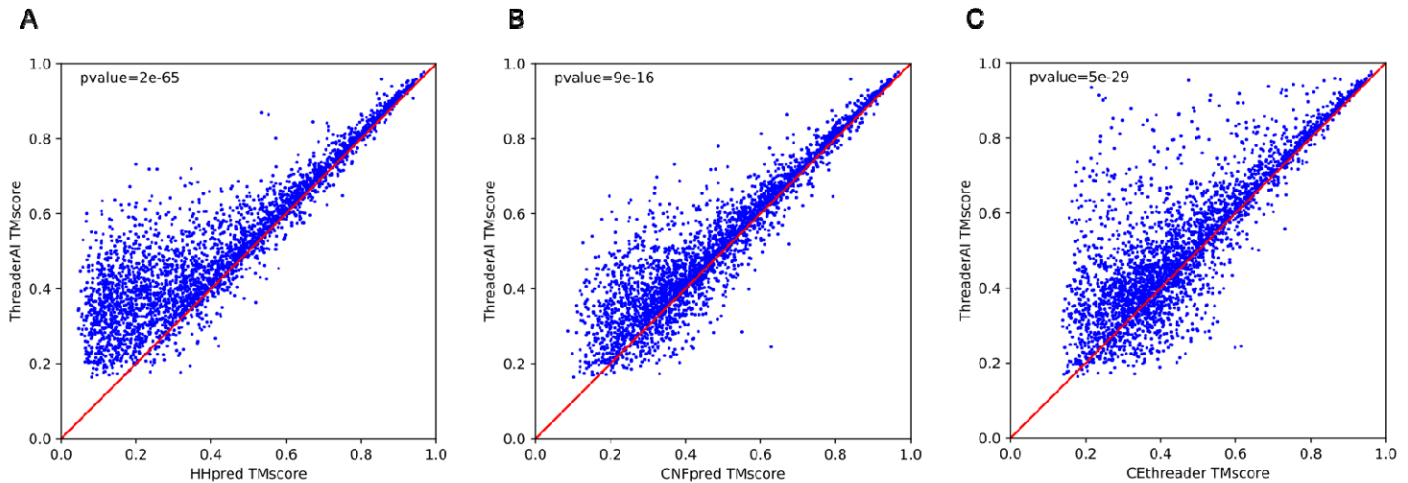
### 270 **3.1 Alignment accuracy on SCOPe3K data**

271 Based on SCOPe's hierarchical classification for proteins, we split all the template-query pairs  
272 into three groups: the pairs similar at family level, at superfamily level, and at fold level. Two  
273 proteins are similar at fold level if both query and template belong to the same fold but different  
274 super families. The similarity at superfamily level and family level are defined in the same way.  
275 Two proteins similar at fold level are conserved in structure but diverges in sequence, and are  
276 usually considered as remote homologs, while two protein similar at family level share high  
277 sequence similarity and are usually considered as close homologs.

278

279 As shown in Table 2 and Figure 2, on SCOPe3K data, ThreaderAI outperforms all other  
280 competitors including HHpred, CNFpred, and CETHreader in terms of alignment accuracy. In  
281 particular, ThreaderAI achieved average TM-score and GDT of 0.510 and 0.437, respectively. In  
282 terms of TM-score, ThreaderAI outperforms HHpred, CNFpred, and CETHreader by 23%, 9%,  
283 and 12%, respectively. The advantage of ThreaderAI over the second-best method is the largest  
284 when the similarity between template and query falls into fold level, which indicates  
285 ThreaderAI's power in modelling of remote homologs. In particular, at the fold level,  
286 ThreaderAI outperforms HHpred, CNFpred, and CETHreader by 56%, 13%, and 11% in terms of  
287 TM-score, respectively. The advantages of ThreaderAI over other methods decreases at the  
288 family level, which is not surprising since it is easy to align two closely-related proteins. At the

289 superfamily level, ThreaderAI outperforms HHpred, CNFpred, and CEthreader by 18%, 9%, and  
 290 14% in terms of TM-score, respectively.



**Figure 2.** Comparison of ThreaderAI and previously published methods using alignment accuracy on SCOPe3K. Each point in the figure represents alignment accuracy of ThreaderAI versus the other competing method.

291  
 292 We also used a t-test to assess the statistical significance of the comparison results. On 3206  
 293 template-query pairs, in terms of TM-score, the *p*-values between ThreaderAI and HHpred,  
 294 CNFpred, and CEthreader are 2e-65, 9e-16, and 5e-29, respectively. Figure 2 shows more details  
 295 on the difference of alignment accuracy between ThreaderAI and the competing methods. In  
 296 terms of TM-score, ThreaderAI achieved better alignment quality than CNFpred for 2743, 2395,  
 297 and 2343 pairs, while worse for 363, 711, and 763 pairs, respectively. It confirms that  
 298 ThreaderAI can generate better alignments than our competing methods.

299

### 300 3.2 Threading performance on CASP13 data

	TBM-easy		TBM-hard		TBM-All	
	TM-score	GDT	TM-score	GDT	TM-score	GDT
ThreaderAI	0.813/0.831	0.752/0.776	0.663/0.702	0.554/0.608	0.761/0.787	0.684/0.719
HHpred	0.753/0.779	0.704/0.733	0.570/0.629	0.477/0.554	0.690/0.728	0.626/0.672
CNFpred	0.785/0.805	0.727/0.751	0.611/0.659	0.520/0.571	0.726/0.755	0.656/0.689
CEthreader	0.770/0.795	0.715/0.740	0.615/0.665	0.533/0.582	0.717/0.751	0.653/0.686

**Table 3.** Threading performance on 61 CASP13 TBM domains. Each cell shows the average quality of the 3D models built from the first-ranked and the best of top five templates.

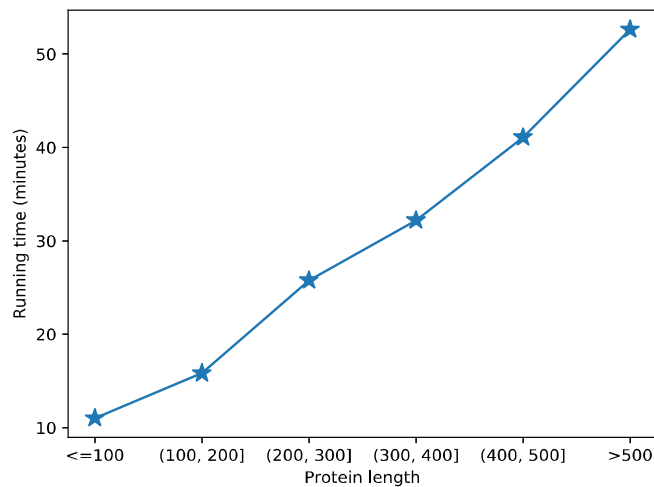
301 We further evaluated the threading performance of our method on the 61 CASP13 TBM domains.  
302 Among the TBM domains, 40 and 21 domains belong to the categories of TBM-easy and TBM-  
303 hard, respectively. Here ThreaderAI and all competitors used the same template database (see  
304 section 2.8.2).

305  
306 As shown in Table 3, on all TBM targets, ThreaderAI outperforms all the competing methods no  
307 matter whether the models are built from the first-ranked or the best of top five templates.

308 ThreaderAI achieves a TM-score 0.761 for first-ranked models, which outperforms HHpred,  
309 CNFpred, and CEthreader 10%, 5%, and 6%, respectively. Overall, ThreaderAI shows larger  
310 advantages on the TBM-hard group in which only remote homologs are available. Specifically,  
311 on TBM-hard group, ThreaderAI outperforms HHpred, CNFpred, and CEthreader by 16%, 9%,  
312 and 8%, respectively. This again indicates ThreaderAI's great advantages in modelling of remote  
313 homologs.

314

### 315 3.3 Running time



**Figure 3.** The running time of ThreaderAI searching query protein in CASP13 data against PDB90. Here we split the data into several groups based on protein lengths. Y-axis is the mean running time in minutes for each group.

316 With the help of GPUs' computational power, ThreaderAI is very efficient in protein threading  
317 (Figure 3). As far as we know, ThreaderAI is the first template-based modelling method which  
318 can take advantage of GPUs. ThreaderAI first uses 3 GeForce-1080 GPUs to generate the  
319 scoring matrices for all templates in the template library and meanwhile uses 4 CPU cores to

320 maintain the data stream for the model. And then ThreaderAI runs the Maximum Accuracy  
321 Algorithm for all scoring matrices on 1 CPU core.

322

323 The running time of ThreaderAI mainly depends on protein length. The protein threading can be  
324 finished within 20 minutes for proteins with less than 200 amino acids. And it takes ThreaderAI  
325 less than 1 hour to finish protein threading even for the proteins with length larger than 500.

326 ThreaderAI is highly scalable as it can use more GPUs.

327

#### 328 **4 Discussion**

329 We developed ThreaderAI, a new template-based method for predicting protein structure using a  
330 deep residual neural network. We show that Threader outperforms the existing popular TBM  
331 methods including HHpred, CNFpred, and CEthreader, in both alignment accuracy and threading  
332 performance, especially on proteins that only have remote homologs with known structure. In  
333 particular, ThreaderAI outperforms CNFpred, another neural network based-method, in which  
334 only one dense layer is used. This demonstrates that advanced neural network models are more  
335 capable of capturing complex sequence-structure relationship.

336

337 ThreaderAI formulates the template-query alignment problem as the classical pixel classification  
338 problem in computer vision. To fulfill this, residue-residue pair scoring is separated from  
339 alignment generation. It's still possible to design an end-to-end model to produce template-query  
340 alignment by combining a deep residual neural network and a chain graphical model such as  
341 Hidden Markov Model<sup>30</sup> and Condition Random Fields<sup>31</sup>. However, in the hybrid model, the  
342 gradients of neural network will entangle with the gradients of chain graphical model which  
343 makes it very inefficient to train a deep model on a large scale of training samples<sup>32</sup>.

344

345 ThreaderAI could be improved in several directions. First, besides deep residual neural network,  
346 other deep learning models such as deep autoregressive models<sup>33</sup> may improve alignment  
347 accuracy. Second, deep attention model<sup>34</sup> may provide a more efficient way to integrate residue-  
348 residue contact information. ThreaderAI integrates residue-residue contacts indirectly by  
349 including the eigenvectors of the contact matrix in which the sign of eigenvectors are decided

350 very heuristically. Local potentials and pairwise potentials related to the residue-residue contact  
351 pairs and non-contacting pairs can be weighted directly with the help of attention mechanisms.

352

353

## 354 **References**

355

- 356 1 Yang, J. Y. *et al.* Improved protein structure prediction using predicted interresidue  
357 orientations. *P Natl Acad Sci USA* **117**, 1496-1503, doi:10.1073/pnas.1914677117 (2020).
- 358 2 Senior, A. W. *et al.* Improved protein structure prediction using potentials from deep  
359 learning. *Nature* **577**, 706-+, doi:10.1038/s41586-019-1923-7 (2020).
- 360 3 Xu, J. B. Distance-based protein folding powered by deep learning. *P Natl Acad Sci USA*  
361 **116**, 16856-16865, doi:10.1073/pnas.1821309116 (2019).
- 362 4 Xu, J. B. & Wang, S. Analysis of distance-based protein structure prediction by deep  
363 learning in CASP13. *Proteins* **87**, 1069-1081 (2019).
- 364 5 Zhu, J. W., Wang, S., Bu, D. B. & Xu, J. B. Protein threading using residue co-variation and  
365 deep learning. *Bioinformatics* **34**, 263-273, doi:10.1093/bioinformatics/bty278 (2018).
- 366 6 Zheng, W. *et al.* Detecting distant-homology protein structures by aligning deep neural-  
367 network based contact maps. *Plos Comput Biol* **15**, doi:ARTN e1007411  
368 10.1371/journal.pcbi.1007411 (2019).
- 369 7 Croll, T. I., Sammito, M. D., Kryshchuk, A. & Read, R. J. Evaluation of template-based  
370 modeling in CASP13. *Proteins* **87**, 1113-1127, doi:10.1002/prot.25800 (2019).
- 371 8 Söding, J. Protein homology detection by HMM–HMM comparison. *Bioinformatics* **21**,  
372 951-960 (2005).
- 373 9 Ma, J. Z., Peng, J., Wang, S. & Xu, J. B. A conditional neural fields model for protein  
374 threading. *Bioinformatics* **28**, i59-i66, doi:10.1093/bioinformatics/bts213 (2012).
- 375 10 Yang, Y. D., Faraggi, E., Zhao, H. Y. & Zhou, Y. Q. Improving protein fold recognition and  
376 template-based modeling by employing probabilistic-based matching between  
377 predicted one-dimensional structural properties of query and corresponding native  
378 properties of templates. *Bioinformatics* **27**, 2076-2082,  
379 doi:10.1093/bioinformatics/btr350 (2011).
- 380 11 Buchan, D. W. A. & Jones, D. T. Eigen THREADER: analogous protein fold recognition by  
381 efficient contact map threading. *Bioinformatics* **33**, 2684-2690,  
382 doi:10.1093/bioinformatics/btx217 (2017).
- 383 12 Peng, J. & Xu, J. B. Boosting Protein Threading Accuracy. *Research in Computational*  
384 *Molecular Biology, Proceedings* **5541**, 31-+ (2009).
- 385 13 Webb, B. & Sali, A. Comparative protein structure modeling using MODELLER. *Current*  
386 *protocols in bioinformatics* **54**, 5.6. 1-5.6. 37 (2016).
- 387 14 Remmert, M., Biegert, A., Hauser, A. & Soding, J. HHblits: lightning-fast iterative protein  
388 sequence searching by HMM-HMM alignment. *Nat Methods* **9**, 173-175 (2012).
- 389 15 Kabsch, W. & Sander, C. Dictionary of Protein Secondary Structure - Pattern-Recognition  
390 of Hydrogen-Bonded and Geometrical Features. *Biopolymers* **22**, 2577-2637 (1983).



- 391 16 Klausen, M. S. *et al.* NetSurfP-2.0: Improved prediction of protein structural features by  
392 integrated deep learning. *Proteins* **87**, 520-527 (2019).
- 393 17 Li, Y., Hu, J., Zhang, C. X., Yu, D. J. & Zhang, Y. ResPRE: high-accuracy protein contact  
394 prediction by coupling precision matrix with deep residual neural networks.  
395 *Bioinformatics* **35**, 4647-4655 (2019).
- 396 18 He, K., Zhang, X., Ren, S. & Sun, J. in *Proceedings of the IEEE conference on computer  
397 vision and pattern recognition*. 770-778.
- 398 19 Clevert, D.-A., Unterthiner, T. & Hochreiter, S. Fast and accurate deep network learning  
399 by exponential linear units (elus). *arXiv preprint arXiv:1511.07289* (2015).
- 400 20 Wang, S., Ma, J., Peng, J. & Xu, J. Protein structure alignment beyond spatial proximity.  
401 *Scientific reports* **3**, 1448 (2013).
- 402 21 Zhang, Y. & Skolnick, J. TM-align: a protein structure alignment algorithm based on the  
403 TM-score. *Nucleic Acids Res* **33**, 2302-2309 (2005).
- 404 22 Loshchilov, I. & Hutter, F. Fixing weight decay regularization in adam. (2018).
- 405 23 Abadi, M. *et al.* TensorFlow: Large-Scale Machine Learning on Heterogeneous  
406 Distributed Systems. *arXiv e-prints* (2016).  
407 <https://ui.adsabs.harvard.edu/abs/2016arXiv160304467A>.
- 408 24 Durbin, R., Eddy, S. R., Krogh, A. & Mitchison, G. *Biological sequence analysis:  
409 probabilistic models of proteins and nucleic acids*. (Cambridge university press, 1998).
- 410 25 Long, J., Shelhamer, E. & Darrell, T. in *Proceedings of the IEEE conference on computer  
411 vision and pattern recognition*. 3431-3440.
- 412 26 Fox, N. K., Brenner, S. E. & Chandonia, J.-M. SCOPe: Structural Classification of  
413 Proteins—extended, integrating SCOP and ASTRAL data and classification of new  
414 structures. *Nucleic Acids Res* **42**, D304-D309 (2014).
- 415 27 Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for  
416 the analysis of massive data sets. *Nature biotechnology* **35**, 1026-1028 (2017).
- 417 28 Zhang, Y. & Skolnick, J. Scoring function for automated assessment of protein structure  
418 template quality. *Proteins: Structure, Function, and Bioinformatics* **57**, 702-710 (2004).
- 419 29 Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database  
420 search programs. *Nucleic Acids Res* **25**, 3389-3402 (1997).
- 421 30 Sisson, S. Hidden Markov models for bioinformatics. *J Roy Stat Soc a Sta* **167**, 194-195,  
422 doi:DOI 10.1111/j.1467-985X.2004.298\_13.x (2004).
- 423 31 Lafferty, J., McCallum, A. & Pereira, F. C. Conditional random fields: Probabilistic models  
424 for segmenting and labeling sequence data. (2001).
- 425 32 Johnson, M. J., Duvenaud, D., Wiltchko, A. B., Datta, S. R. & Adams, R. P. Composing  
426 graphical models with neural networks for structured representations and fast inference.  
427 *Adv Neur In* **29** (2016).
- 428 33 Yang, Z. *et al.* in *Advances in neural information processing systems*. 5754-5764.
- 429 34 Vaswani, A. *et al.* in *Advances in neural information processing systems*. 5998-6008.
- 430  
431