

# **Predicting changes in protein thermodynamic stability upon point mutation with deep 3D convolutional neural networks**

Bian Li<sup>1,2,4</sup>, Yucheng T. Yang<sup>1,2</sup>, John A. Capra<sup>4\*</sup>, Mark B. Gerstein<sup>1,2,3\*</sup>

<sup>1</sup> Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT 06520, USA

<sup>2</sup> Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT 06520, USA

<sup>3</sup> Department of Computer Science, Yale University, New Haven, CT 06520, USA

<sup>4</sup> Department of Biological Sciences and Vanderbilt Genetics Institute, Vanderbilt University, Nashville, TN 37235, USA

\* Corresponding authors

E-mails: [tony.capra@vanderbilt.edu](mailto:tony.capra@vanderbilt.edu) (J.A.C.) and [pi@gersteinlab.org](mailto:pi@gersteinlab.org) (M.B.G.)

## 26 **Abstract**

27 Predicting mutation-induced changes in protein thermodynamic stability ( $\Delta\Delta G$ ) is of great  
28 interest in protein engineering, variant interpretation, and understanding protein biophysics. We  
29 introduce ThermoNet, a deep, 3D-convolutional neural network designed for structure-based  
30 prediction of  $\Delta\Delta G$ s upon point mutation. To leverage the image-processing power inherent in  
31 convolutional neural networks, we treat protein structures as if they were multi-channel 3D  
32 images. In particular, the inputs to ThermoNet are uniformly constructed as multi-channel voxel  
33 grids based on biophysical properties derived from raw atom coordinates. We train and evaluate  
34 ThermoNet with a curated data set that accounts for protein homology and is balanced with  
35 direct and reverse mutations; this provides a framework for addressing biases that have likely  
36 influenced many previous  $\Delta\Delta G$  prediction methods. ThermoNet demonstrates performance  
37 comparable to the best available methods on the widely used  $S^{\text{sym}}$  test set. However,  
38 ThermoNet accurately predicts the effects of both stabilizing and destabilizing mutations, while  
39 most other methods exhibit a strong bias towards predicting destabilization. We further show  
40 that homology between  $S^{\text{sym}}$  and widely used training sets like S2648 and VariBench has likely  
41 led to overestimated performance in previous studies. Finally, we demonstrate the practical  
42 utility of ThermoNet in predicting the  $\Delta\Delta G$ s for two clinically relevant proteins, p53 and  
43 myoglobin, and for pathogenic and benign missense variants from ClinVar. Overall, our results  
44 suggest that 3D convolutional neural networks can model the complex, non-linear interactions  
45 perturbed by mutations, directly from biophysical properties of atoms.

## 46 **Author Summary**

47 The thermodynamic stability of a protein, usually represented as the Gibbs free energy for the  
48 biophysical process of protein folding ( $\Delta G$ ), is a fundamental thermodynamic quantity. Predicting  
49 mutation-induced changes in protein thermodynamic stability ( $\Delta\Delta G$ ) is of great interest in protein  
50 engineering, variant interpretation, and understanding protein biophysics. However, predicting

51  $\Delta\Delta G$ s in an accurate and unbiased manner has been a long-standing challenge in the field of  
52 computational biology. In this work, we introduce ThermoNet, a deep, 3D-convolutional neural  
53 network designed for structure-based  $\Delta\Delta G$  prediction. To leverage the image-processing power  
54 inherent in convolutional neural networks, we treat protein structures as if they were multi-  
55 channel 3D images. ThermoNet demonstrates performance comparable to the best available  
56 methods. However, ThermoNet accurately predicts the effects of both stabilizing and  
57 destabilizing mutations, while most other methods exhibit a strong bias towards predicting  
58 destabilization. We also demonstrate that the presence of homologous proteins in commonly  
59 used training and testing sets for  $\Delta\Delta G$  prediction methods has likely influenced previous  
60 performance estimates. Finally, we highlight the practical utility of ThermoNet by applying it to  
61 predicting the  $\Delta\Delta G$ s for two clinically relevant proteins, p53 and myoglobin, and for pathogenic  
62 and benign missense variants from ClinVar.

## 63 **Introduction**

64 The thermodynamic stability of a protein, usually represented as the Gibbs free energy for the  
65 biophysical process of protein folding ( $\Delta G$ ), is a fundamental thermodynamic quantity. The  
66 magnitude of  $\Delta G$  is collectively determined by the intramolecular interactions between amino  
67 acid residues within the protein and the interactions between the protein and the physiological  
68 environment around it (1). When a mutation causes amino acid substitution in a protein, it is  
69 likely that the stability of the mutant protein will be affected compared to the wild type. (Note that  
70 the term “wild type” is not preferred because humans have substantial protein-coding genetic  
71 diversity (2). A better term would be a “reference state”. However, as it will not cause confusion  
72 in this work and for consistency with previous work, we use “wild type” throughout the text.) This  
73 change in protein thermodynamic stability (i.e.  $\Delta\Delta G$ ) incurred by mutation is of fundamental  
74 importance to medicine and biotechnology. Many disease-causing mutations are single-point  
75 amino acid substitutions that lead to a substantial  $\Delta\Delta G$  of the corresponding protein, and such

76 single-point mutations are a key mechanism underlying a wide spectrum of molecular disorders  
77 (3-5). Given the huge number of variants discovered by large-scale population-level exome and  
78 genome sequencing studies and clinical genetic tests, there is a tremendous interest in  
79 predicting whether these variants are likely to exert any impact on protein function. In addition,  
80 in developing new biopharmaceuticals, one of the early goals is usually to design proteins with  
81 the intended thermodynamic stability. However, this task is often laborious, if not impossible,  
82 and usually involves experimentally screening an enormous number of mutant proteins (6).  
83 Thus, it is desirable to have an efficient and accurate computational tool to prioritize the set of  
84 mutant proteins to be experimentally tested.

85 Toward these goals, several programs have been developed for estimating  $\Delta\Delta G$ s. These  
86 methods either rely on explicit biophysical modeling of amino acid interactions coupled with  
87 conformational sampling of protein structures (7-11) or apply machine/statistical learning to  
88 extract patterns from various types of amino acid sequence, evolutionary, and/or protein  
89 structural features (12-20). While these methods have been useful in many applications (21, 22),  
90 they have substantial limitations. For example, physics-based methods are computationally  
91 demanding and low-throughput; these challenges have largely prevented them from being  
92 applied to large-scale protein engineering and variant interpretation tasks. On the other hand,  
93 several studies have highlighted significant bias in the predictions of machine learning-based  
94 methods; they tend to predict mutations as destabilizing more often than stabilizing (19, 23-25).  
95 The main source of this bias likely comes from the fact that the training sets are dominated by  
96 experiment-derived destabilizing mutations and that machine learning methods are prone to  
97 overfitting to training sets (24, 26). Thus, there is a need for new methods that can make  
98 quantitative, unbiased prediction of  $\Delta\Delta G$ s with high throughput.

99 Here, we describe ThermoNet, a computational framework based on deep 3D convolutional  
100 neural networks (3D-CNNs) for predicting  $\Delta\Delta G$ s upon single-point mutation. We model the  
101 structure of each mutation assuming that single-point mutations introduce negligible

102 perturbation to the overall architecture of protein structure. We treat protein structures as if they  
103 were 3D images with voxels parameterized using atom biophysical properties (27, 28). We  
104 leverage the power of the architecture of CNNs in detecting spatially proximate features. These  
105 local biochemical interaction detectors are then hierarchically composed into more intricate  
106 features with the potential to describe the complex and nonlinear phenomenon of molecular  
107 interaction. We address the bias in many previous methods towards predicting destabilization  
108 by training ThermoNet on a balanced data set generated through anti-symmetry-based data  
109 augmentation, i.e. for each mutation, we consider both the direct and reverse versions. We  
110 further demonstrate and address an unappreciated source of bias in previous performance  
111 estimates due to homology between training and evaluation sets. We show that ThermoNet  
112 achieves state-of-the-art performance comparable to previously developed methods on a widely  
113 used test set with minimal prediction bias. We also demonstrate the applicability of ThermoNet  
114 by showing that ThermoNet accurately predicts the  $\Delta\Delta G$ s of the missense mutations in two  
115 biologically important proteins, the p53 tumor suppressor protein and myoglobin and that  
116 ThermoNet-predicted  $\Delta\Delta G$ s of ClinVar missense variants fall within the experimentally observed  
117 range and are consistent with the expectations of a biophysical model of protein evolution.

## 118 **Results**

### 119 ***An overview of ThermoNet***

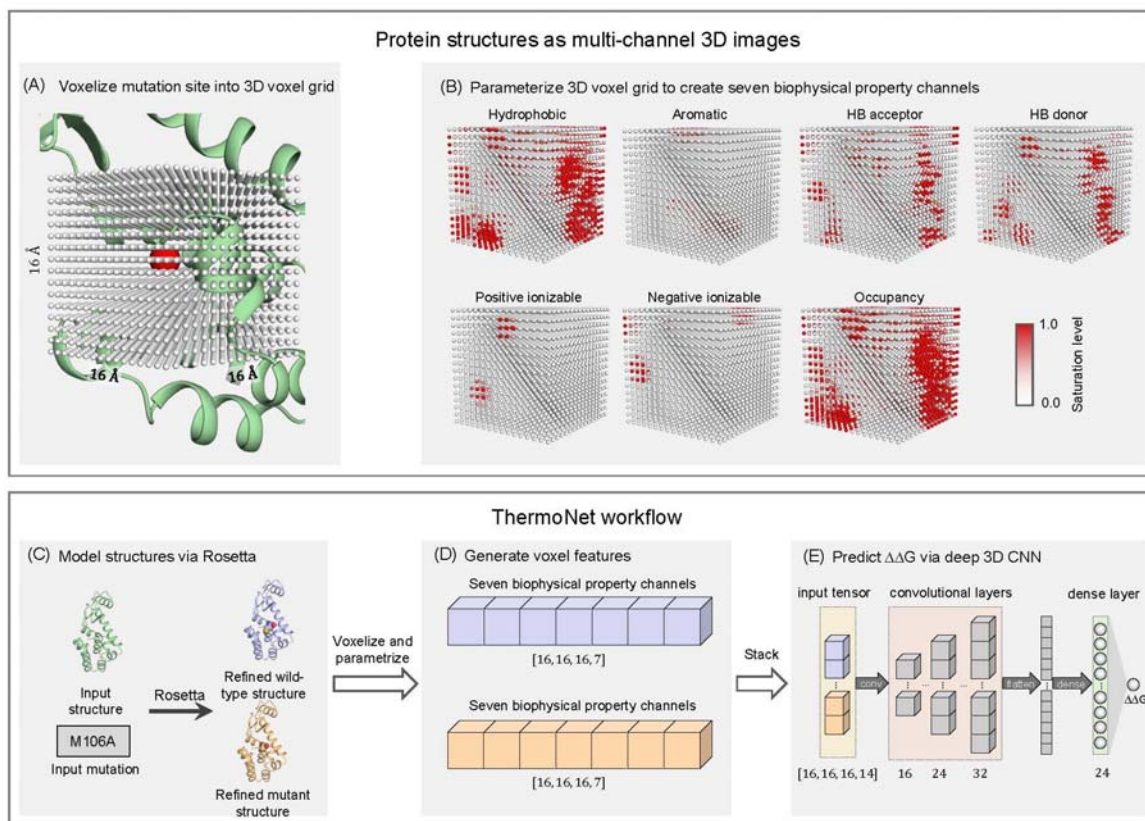
120 To predict the  $\Delta\Delta G$  of a point mutation, we take advantage of recent advances in deep learning  
121 for computer vision (29) and the successes of deep convolutional neural networks in biophysical  
122 problems (27, 28, 30-32). We treat protein structures as if they were 3D images with voxels  
123 parameterized using atom biophysical properties (27, 28) (Fig 1A,1B, and Table 1). For a given  
124 single-point mutation, ThermoNet requires that a 3D structure (either experimentally determined  
125 or modeled via homology modeling) of one of the alleles is available. As a first step, ThermoNet  
126 constructs a structural model for the mutant from the structure of the wild type using the Rosetta

127 macromolecular modeling suite (Methods) (9, 33). ThermoNet assumes that the  $\Delta\Delta G$  of a point  
128 mutation can be sufficiently captured by modeling the 3D physicochemical environment around  
129 the mutation site. It thus extracts predictive features by treating protein structures as if they were  
130 3D images and voxelizing the space around the mutation site of both the wild-type structure and  
131 the corresponding mutant structural model (Fig 1C). Each voxel is parameterized with seven  
132 predefined rules (Table 1) to characterize the physicochemical nature of its neighboring atoms.  
133 The feature maps are then stacked to create a tensor with size [16,16,16,14] as input to the  
134 trained ensemble of ten deep 3D-CNNs, which generate a prediction of the  $\Delta\Delta G$  the given  
135 mutation causes to the wild-type structure (Fig 1D, 1E, and Methods). Each of the component  
136 3D-CNN models consists of three 3D convolutional layers with 16, 24, and 32 neurons  
137 respectively and one densely connected layer of 24 neurons (Fig 1E). These architectural  
138 hyperparameters (i.e. number of neurons in convolutional and densely connected layers and  
139 sizes of the input voxel grid) were tuned via five-fold cross-validation (Methods, Fig S1).

140 Table 1. Chemical property channels (AutoDock4 atom types) of a protein structure voxel

Property	Rule
Hydrophobic	Aliphatic or aromatic carbon atoms
Aromatic	Aromatic carbon atoms
Hydrogen bond donor	Nitrogen, oxygen, sulfur atoms with lone-pair electrons
Hydrogen bond acceptor	Polar hydrogen atoms
Positive ionizable	Atoms with positive charge
Negative ionizable	Atoms with negative charge
Occupancy	All atom types

141



142

143

Fig 1. An overview of the ThermoNet computational framework.

144 (A) Protein structures are treated as if they were 3D images. A  $16 \text{ \AA} \times 16 \text{ \AA} \times 16 \text{ \AA}$  cubic  
145 neighborhood centered at the  $C_{\beta}$  atom (red sphere) of the mutated residue (or  $C_{\alpha}$  atom in the  
146 case of a glycine) of an example protein (PDB ID: 1L63) is discretized into a 3D voxel grid at a  
147 resolution of  $1 \text{ \AA}$ . Each voxel is represented by a gray dot. (B) Just as an RGB image has three  
148 color channels, the 3D voxel grid is parameterized with seven chemical property channels:  
149 hydrophobic, aromatic, hydrogen bonding donor, hydrogen bond acceptor, positive ionizable,  
150 negative ionizable, and occupancy. The saturation level of each voxel ranges from 0.0 to 1.0  
151 and is colored accordingly (Methods). (C) To predict the change in thermodynamic stability  
152 caused by a given single-point mutation, ThermoNet calls Rosetta to refine the wild-type  
153 structure and to create a structural model of the mutant protein. (D) ThermoNet voxelizes the  
154 space around the mutation site of both the Rosetta-refined wild-type structure and the  
155 corresponding mutant structural model. Both the 3D voxel grid of the wild-type structure and that  
156 of the mutant model are parameterized accordingly to create two  $[16, 16, 16, 7]$  feature maps. (E)  
157 The feature maps are then stacked to create a  $[16, 16, 16, 14]$  tensor as an input to the trained  
158 deep 3D convolutional neural network. The final output of the network is the predicted  $\Delta\Delta G$  the  
159 given mutation causes to the wild-type protein structure.

## 160 ***Creating data sets for robust training and testing of ThermoNet***

161 The ability of a machine-learning model to generalize can be overestimated when there is data  
162 leakage between the training set and the test set. In structural bioinformatics problems such as  
163  $\Delta\Delta G$  prediction, such data leakage can result when the training set contains proteins that are  
164 homologous to proteins in the test set. This is because the effects of different mutations in the  
165 same protein or homologous proteins are not necessarily independent. In most previous  
166 methods for  $\Delta\Delta G$  prediction, this data leakage issue was not fully appreciated, and homologous  
167 proteins were present between training and test sets. For example, a recent method used  
168 randomly selected subsets of mutants from the widely used S2648 data set for training and  
169 testing (34). Not surprisingly, the training and test sets of shared 61 identical proteins  
170 (Supporting Information Table S1).

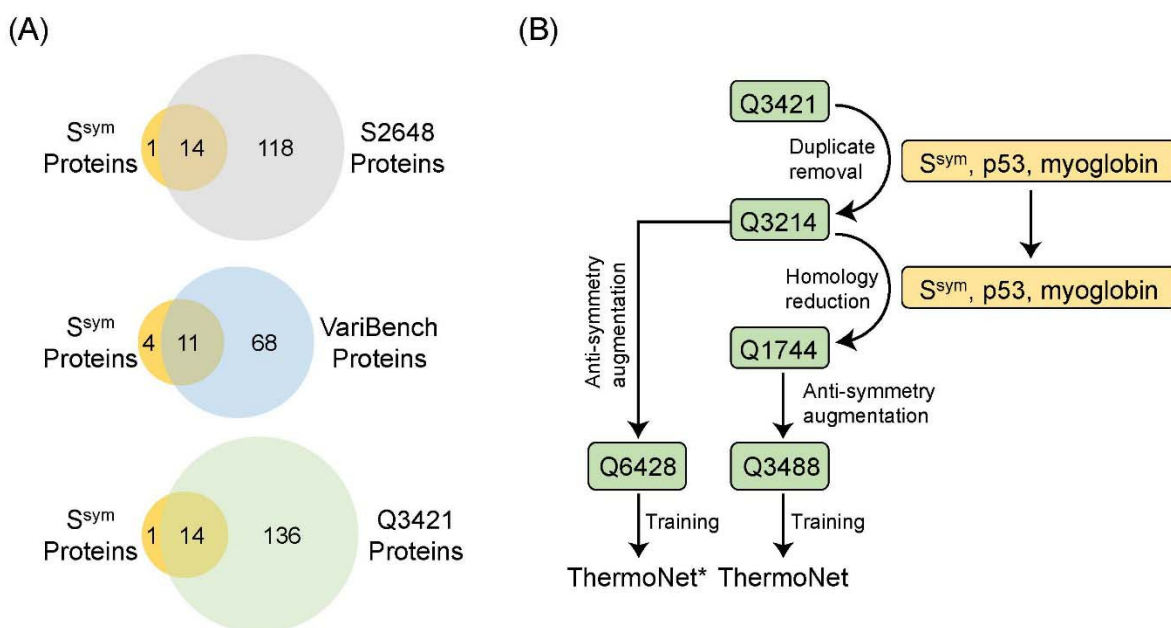
171 The issue of having mutations from the same protein in both training and test sets was  
172 addressed in developing the mCSM and INPS methods, where the cross-validation procedure  
173 ensured that mutations of the same protein remained together in either the training or test set  
174 (16, 19). While this was a step in the right direction, grouping mutations at the protein level is  
175 not sufficient to remove the homology between training and test sets. For example, we found  
176 substantial homology between 132 proteins within S2648 (Supporting Information Fig S2A).  
177 Thus, splitting the S2648 data set for training and testing at the protein level is likely to end up  
178 with shared homology between the splits. In fact, using the PISCES server (35) to remove  
179 redundancy from the S2648 data set resulted in only 104 non-redundant proteins (out of 132) at  
180 the level of  $< 25\%$  sequence identity (Supporting Information Table S2). Homology is common  
181 among data sets used in training  $\Delta\Delta G$  predictors; for example, many proteins in the VariBench  
182 (36) data set also share substantial homology (Fig S2B).

183 We further highlight that the widely used S2648 data set includes fourteen proteins from  $S^{\text{sym}}$   
184 data set, which was previously used to evaluate a wide range of  $\Delta\Delta G$  predictors (24, 37), and



185 another eight proteins that are putative homologs to proteins in  $S^{\text{sym}}$  (Fig 2A, Supporting  
186 Information Table S3). In addition, a similar level of overlap also exists between the VariBench  
187 and Q3421 data sets and  $S^{\text{sym}}$  (Fig 2A Supporting Information Table S4 and S5). Thus,  
188 performance estimates on  $S^{\text{sym}}$  of methods trained using S2648 or parameterized using  
189 VariBench are likely to be overly optimistic. For example, in a recent publication, a version of  
190 INPS trained using a data set obtained by removing all proteins with > 25% sequence identity to  
191 proteins in  $S^{\text{sym}}$  showed substantially reduced performance (38).

192 Thus, to train and evaluate ThermoNet, we implemented a rigorous procedure to reduce the  
193 sequence similarity between the training and test sets ( $S^{\text{sym}}$ , p53, myoglobin) by removing  
194 duplicate data points and pruning protein-level homology (Fig 2B). Our rigorous pruning of the  
195 starting Q3421 data set (15) resulted in a data set consisting of 1,744 distinct mutations. This  
196 data set was then augmented by creating a reverse mutation data point for each of the 1,744  
197 direct mutations according to the anti-symmetry property of  $\Delta\Delta G$  (Methods), thus giving to a  
198 total of 3,488 data points for the training of ThermoNet (Fig 2B). While the pruning nearly halved  
199 the size of available training data, as discussed in the following section, models trained using  
200 the resulting augmented Q3488 data set will be less likely to have an overestimated  
201 performance when evaluated on  $S^{\text{sym}}$ . Finally, to explore the influence of homology between  
202 training and test sets on estimates of model performance, we also augmented Q3214, the  
203 intermediate data set before the step of homology reduction, to train a different version of  
204 ThermoNet, called ThermoNet\* (Fig 2B).



205  
 206 Fig 2. Data set curation and identification of shared homology.  
 207 (A) Venn diagrams showing the amount of overlap at the protein level between three widely  
 208 used training sets S2648, VariBench, and Q3421 for  $\Delta\Delta G$  predictors and the  $S^{\text{sym}}$  test set.  
 209 Numbers in these diagrams indicate protein counts. Upper panel and lower panel indicate that  
 210 both S2648 and Q3421 share 14 identical proteins with  $S^{\text{sym}}$ ; middle panel indicates that  
 211 VariBench and  $S^{\text{sym}}$  share 11 identical proteins. All three data sets share additional homology  
 212 with  $S^{\text{sym}}$ , which is presented in Supporting Information Table S3, S4, and S5, respectively. (B)  
 213 Creating data sets for robust training and testing of ThermoNet. We started with the Q3421 set  
 214 of 3421 mutations from 150 proteins. (Numbers in data set names indicate the number of  
 215 unique mutations the data set contains.) After homology reduction and anti-symmetry data  
 216 augmentation (Methods), this data curation workflow gives a training set of 3488 mutations with  
 217 an equal representation of stabilizing and destabilizing changes and reduced homology to the  
 218  $S^{\text{sym}}$  test set. A separate data set called Q6428 was also created by augmenting the Q3214 data  
 219 set before homology reduction to train ThermoNet\*.

## 220 ***ThermoNet achieves state-of-the-art performance on blind test set***

221 We systematically compared ThermoNet and ThermoNet\* with seventeen  $\Delta\Delta G$  predictors on  
 222 the  $S^{\text{sym}}$  balanced data set to evaluate their performance and degree of bias with respect to the  
 223  $\Delta\Delta G$  anti-symmetry between direct and reverse mutations (Methods). A brief summary of the  
 224 characteristics of these  $\Delta\Delta G$  predictors and their references are given in Supporting Information  
 225 Table S6. In short, the predictors are based on diverse features and strategies with some, like  
 226 ThermoNet based only on structural information, while others like DDGun3D (37) and STRUM

227 (15), integrate structural information with sequence and evolutionary features. The  $S^{\text{sym}}$  data set,  
228 which was constructed previously for assessing the biases of  $\Delta\Delta G$  predictors (24). It consists of  
229 experimentally measured  $\Delta\Delta G$  values for 342 direct and the corresponding reverse mutations (a  
230 total of 684 mutations) from fifteen protein chains for which the structures of both the wild-type  
231 and mutant proteins have been resolved by X-ray crystallography with a resolution of 2.5 Å or  
232 better. This data set is by construction balanced with respect to stabilizing and destabilizing  
233 mutations, thus enabling the evaluation of prediction bias. However, as noted in the previous  
234 section, many proteins in  $S^{\text{sym}}$  overlap or are homologous to proteins in commonly used training  
235 sets (Fig 2A).

236 To evaluate performance, we computed the root mean square error  $\sigma$  and the Pearson  
237 correlation coefficient  $r$  separately for direct and reverse mutations. We measured prediction  
238 bias by two statistics, the Pearson correlation coefficient  $r_{dir-rev}$  between the predictions for  
239 direct and those for reverse mutations and the  $\delta$  value, defined as:  $\delta = \Delta\Delta G_{rev} + \Delta\Delta G_{dir}$ . A  
240 perfectly unbiased predictor would give  $r_{dir-rev} = -1$  and  $\langle\delta\rangle = 0 \text{ kcal/mol}$ .

241 ThermoNet achieves strong prediction accuracy that is comparable for direct mutations  
242 ( $r_{dir} = 0.47$  and  $\sigma_{dir} = 1.56 \text{ kcal/mol}$ ; Table 2, Fig 3A) and the corresponding set of reverse  
243 mutations ( $r_{rev} = 0.47$  and  $\sigma_{rev} = 1.55 \text{ kcal/mol}$ , Fig 3C). This suggests that ThermoNet did not  
244 overfit to direct mutations. The fractions of mutations for which the prediction error is within  
245  $0.5 \text{ kcal/mol}$  and  $1.0 \text{ kcal/mol}$  are 36.3% and 58.8% for direct mutations and 36.5% and 58.5%  
246 for reverse mutations (Fig 3B and 3D). ThermoNet successfully reduces prediction bias with a  
247 near-perfect  $r_{dir-rev}$  (-0.96) and a negligible  $\langle\delta\rangle$  (-0.01) (Fig 3E). We also report the distribution  
248 of  $\delta$ , since  $\langle\delta\rangle$  cannot distinguish large, but symmetric, bias from low bias (Methods). As shown  
249 in Fig 3F, 40.9% and 96.2% of mutations have a prediction bias  $< 0.1 \text{ kcal/mol}$  and  $<$   
250  $0.5 \text{ kcal/mol}$ , respectively.

251 Table 2. Comparative analysis using the balanced blind test set  $S^{\text{sym}}$ .

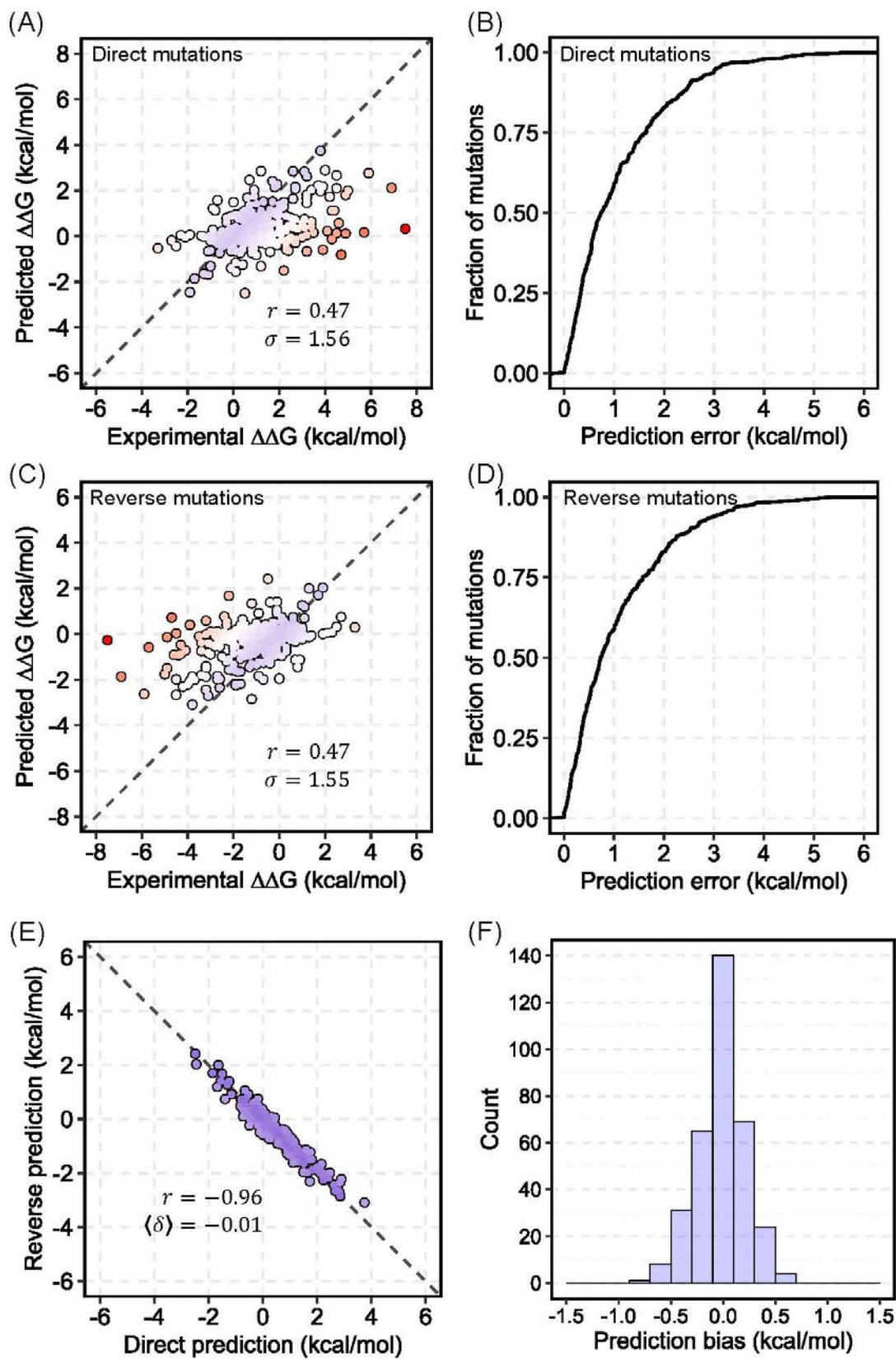
Method	$\sigma_{dir}$	$r_{dir}$	$\sigma_{rev}$	$r_{rev}$	$r_{dir-rev}$	$\langle\delta\rangle$
ThermoNet*	1.42	0.58	1.38	0.59	-0.95	-0.05
DDGun3D	1.42	0.56	1.46	0.53	-0.99	-0.02
DDGun	1.47	0.48	1.50	0.48	-0.99	-0.01
ThermoNet	1.56	0.47	1.55	0.47	-0.96	-0.01
PoPMuSiC <sup>sym</sup>	1.58	0.48	1.62	0.48	-0.77	0.03
MAESTRO	1.36	0.52	2.09	0.32	-0.34	-0.58
FoldX	1.56	0.63	2.13	0.39	-0.38	-0.47
PoPMuSiC 2.1	1.21	0.63	2.18	0.25	-0.29	-0.71
SDM	1.74	0.51	2.28	0.32	-0.75	-0.32
iSTABLE	1.10	0.72	2.28	-0.08	-0.05	-0.60
I-Mutant 3.0	1.23	0.62	2.32	-0.04	0.02	-0.68
NeEMO	1.08	0.72	2.35	0.02	0.09	-0.60
DUET	1.20	0.63	2.38	0.13	-0.21	-0.84
mCSM	1.23	0.61	2.43	0.14	-0.26	-0.91
MUPRO	0.94	0.79	2.51	0.07	-0.02	-0.97
STRUM	1.05	0.75	2.51	-0.15	0.34	-0.87
Rosetta	2.31	0.69	2.61	0.43	-0.41	-0.69
AUTOMUTE	1.07	0.73	2.61	-0.01	-0.06	-0.99
CUPSAT	1.71	0.39	2.88	0.05	-0.54	-0.72

252  $\sigma_{dir}$  and  $r_{dir}$  are the root mean square deviation and the Pearson correlation coefficient between the  
 253 predicted and experimental  $\Delta\Delta G$  values for the direct mutations in  $S^{\text{sym}}$ . Many of these mutations belong  
 254 to the training set of the machine-learning-based methods tested (24), so their performances are likely to  
 255 be overestimated.  $\sigma_{inv}$  and  $r_{rev}$  are the root mean square deviation and the Pearson correlation coefficient  
 256 between the predicted and experimental  $\Delta\Delta G$  values for the reverse mutations in  $S^{\text{sym}}$ . These mutations  
 257 do not belong to the training data sets and thus constitute an independent test set. The parameter  $\delta$   
 258 quantifies the prediction bias and is defined as:  $\delta = \Delta\Delta G_{rev} + \Delta\Delta G_{dir}$ . A perfectly non-biased tool should  
 259 have  $\delta = 0$  for every mutation. We used here its average value  $\langle\delta\rangle$  taken over all mutations that belong to  
 260  $S^{\text{sym}}$ . Numbers for DDGun and DDGun3D were obtained from reference (37) and those for all other  
 261 methods were obtained from (24). The methods are ranked according to their performance,  
 262  $\sigma_{rev}$ , on reverse mutations. Both ThermoNet\* and ThermoNet were trained using a data set balanced with  
 263 direct and reverse mutations, but the data set for training ThermoNet was not homology-reduced with  
 264 respect to  $S^{\text{sym}}$  (Fig 2B).  
 265

266 We also report the performance of ThermoNet\*, different version of ThermoNet trained using  
 267 the Q6428 data set augmented from the intermediate data set Q3214 before the step of  
 268 homology reduction in our data set curation procedure (Fig 2B and Methods). ThermoNet\* was  
 269 trained in the exact same way as ThermoNet except that the homology of its training set Q3214  
 270 to  $S^{\text{sym}}$  was retained. Thus, the parameterization of ThermoNet\* is comparable to previous  
 271 methods that did not consider homology. As expected, evaluation of ThermoNet\* on  $S^{\text{sym}}$  shows

272 even better performance in  $\sigma_{rev}$  (1.38 vs. 1.55 kcal/mol) and  $r_{rev}$  (0.59 vs. 0.48) than  
273 ThermoNet, which was trained using the data set obtained after homology reduction (Table 2,  
274 Fig 2B, and Methods). Thus, the performance of many previously developed methods is likely to  
275 be substantially lower if they had been trained using a data set that shared no homology with  
276  $S^{sym}$ . In contrast, ThermoNet's strong performance, even after removing homology reduction,  
277 suggests robust generalization in real-life applications.

278 Compared to other  $\Delta\Delta G$  predictors, ThermoNet\* achieves the best performance on reverse  
279 mutations, and the methods that outperform it on direct mutations all have substantial bias  
280 against reverse mutations ( $\sigma_{rev} > 2.09$  kcal/mol and  $\langle \delta \rangle < -0.58$ ). The seemingly good  
281 performance of many machine learning-based methods on direct mutations, but poor  
282 performance on reverse mutations suggests potential overfitting due to unbalanced training sets  
283 (24-26, 39). ThermoNet also performs well, but as a result of the reduction in performance due  
284 to removing homology between training and validation sets, the DDGun and DDGun3D methods  
285 outperform it on direct and reverse mutations. Unfortunately, it is not possible to retrain and  
286 evaluate all the other methods on the homology pruned training set, so we cannot directly  
287 compare the other methods to ThermoNet. Nonetheless, the fact that it still outperforms most  
288 suggests its utility and robustness.



290 Fig 3. Performance of ThermoNet on the blind test set.

291 (A) Performance of ThermoNet on predicting  $\Delta\Delta G$  for direct mutations; The Pearson correlation  
292 coefficient ( $r$ ) between predicted values and experimentally determined values is 0.47, and the  
293 root-mean-square deviation ( $\sigma$ ) of predicted values from experimentally determined values is  
294 1.56 kcal/mol. The dots are colored in gradient from blue to red such that blue represents the  
295 most accurate prediction and red indicates the least accurate prediction. (B) Cumulative  
296 distribution of ThermoNet prediction error on direct mutations. (C) Performance of ThermoNet  
297 on predicting  $\Delta\Delta G$  for the reverse mutations ( $r = 0.47$ ,  $\sigma = 1.55$  kcal/mol). (D) Cumulative  
298 distribution of ThermoNet prediction error on reverse mutations. (E) Direct versus reverse  $\Delta\Delta G$   
299 values of all the mutations in the blind test set predicted by ThermoNet. A perfectly unbiased  
300 predictor would give  $r = -1$  and  $\langle\delta\rangle = 0$  kcal/mol. ThermoNet successfully reduces prediction  
301 bias with  $r = -0.96$  and  $\langle\delta\rangle = -0.01$  kcal/mol. (F) Distribution of ThermoNet prediction bias.

### 302 ***Structural models of reverse mutations are necessary for unbiased $\Delta\Delta G$ predictions***

303 To evaluate whether the inclusion of the reverse mutations is necessary for the reduction in  
304 prediction bias, we trained a predictor following the same procedure for training ThermoNet but  
305 using a data set consisting of only the 1,744 direct mutations and their associated experimental  
306  $\Delta\Delta G$ s (i.e. the Q1744 data set in Fig 2B). We applied this predictor to predict the  $\Delta\Delta G$ s of the  
307 direct and reverse mutations of the  $S^{\text{sym}}$  test set. As shown in Fig S3,  $\Delta\Delta G$ s of direct mutations  
308 predicted by these models correlate reasonably well ( $r = 0.47$  and  $\sigma = 1.38$  kcal/mol) with the  
309 experimental values and are comparable to the performance of the ensemble of networks  
310 trained using the balanced data set Q3488. In contrast, these models perform poorly ( $r = -0.06$   
311 and  $\sigma = 2.40$  kcal/mol) in predicting the  $\Delta\Delta G$ s of the corresponding set of reverse mutations (Fig  
312 S3). This suggests that the models were biased toward the training set which is dominated by  
313 destabilizing mutations. This is confirmed by the strongly positive correlation between the  
314 predictions for direct mutations and those for reverse mutations and the large prediction bias  
315 ( $r_{dir-rev} = 0.35$  and  $\langle\delta\rangle = 1.63$  kcal/mol) (Fig S3). Compared to the performance of ThermoNet,  
316 which was trained using the balanced data set Q3488, the results highlight the necessity of a  
317 balanced data set for correcting prediction bias.

318 **Case studies: The p53 tumor suppressor protein and myoglobin.**

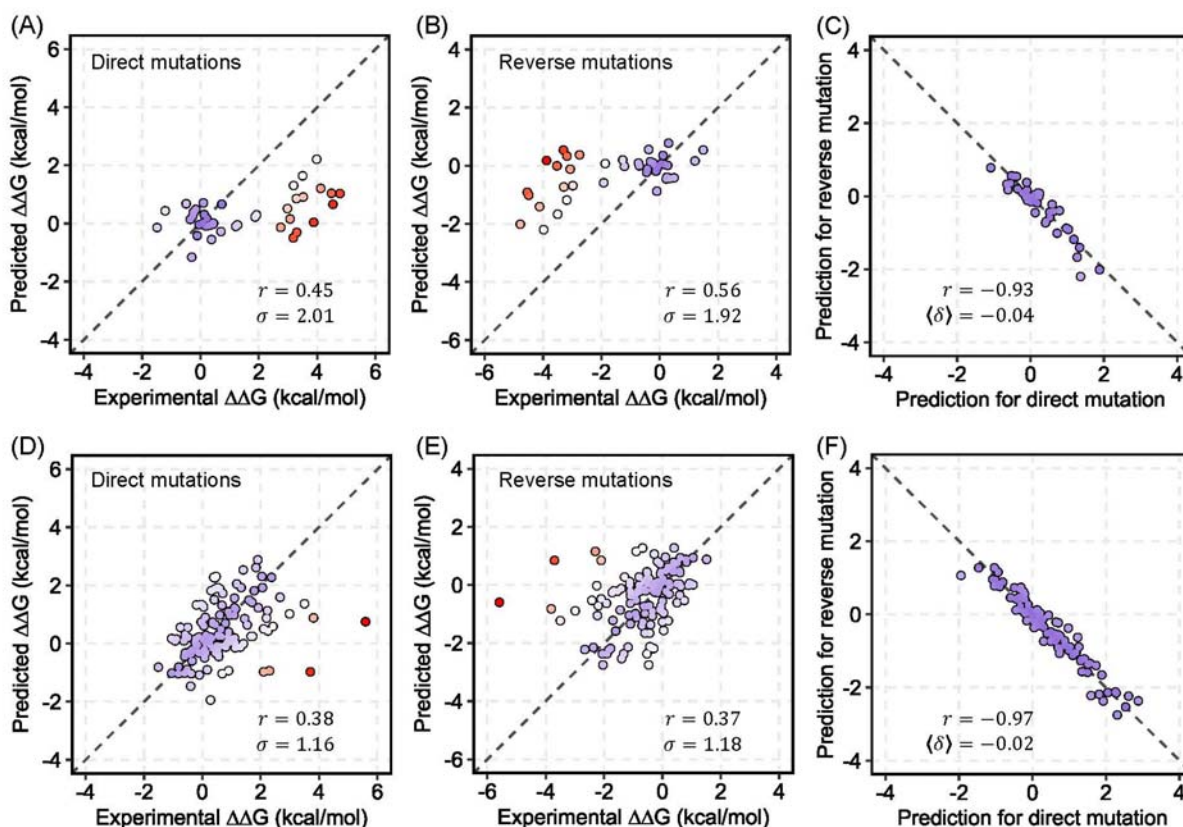
319 We further tested ThermoNet by predicting the  $\Delta\Delta G$ s of single-point mutations in the p53 tumor  
320 suppressor protein and myoglobin whose thermodynamic effects have previously been  
321 experimentally measured. The *TP53* tumor suppressor gene encodes the p53 transcription  
322 factor that is mutated in ~45% of all human cancers (Beroud and Soussi, 2003; Olivier et al.  
323 2002). Unlike most tumor suppressor proteins that are inactivated by deletion or truncation  
324 mutations, single amino acid substitutions in p53 often modify DNA binding or disrupt the  
325 conformation and stability of p53 (Olivier et al., 2002). Myoglobin is a cytoplasmic globular  
326 protein that regulates cellular oxygen concentration in cardiac myocytes and oxidative skeletal  
327 muscle fibers by reversible binding of oxygen through its heme prosthetic group (40). The p53  
328 data set consists of 42 mutations compiled in a previous study (16) within the DNA binding  
329 domain of p53. The myoglobin data set consists of 134 mutations scattered throughout the  
330 protein chain compiled in a previous study (41). We note that none of the mutations in these two  
331 data sets were present in the training set and that proteins that are likely to be homologous to  
332 p53 and myoglobin were also removed from the training set of ThermoNet (Fig 2B, Methods).  
333 We used published crystal structures of p53 (PDB ID: 2OCJ) and myoglobin (PDB ID: 1BZ6) to  
334 create one structural model for each of the mutations in these two data sets respectively using  
335 the *FastRelax* protocol in Rosetta (42). These predictions were compared directly with the  
336 experimentally determined thermodynamic effects (Fig 4).

337 For p53,  $\Delta\Delta G$ s of both direct and reverse mutations predicted by ThermoNet correlate with the  
338 experimental measurements ( $r = 0.45$  and  $0.56$ ) and have little bias ( $r_{dir-rev} = -0.93$  and  $\langle\delta\rangle =$   
339  $-0.04$  kcal/mol). However, the predicted  $\Delta\Delta G$ s of myoglobin mutations correlate less well  
340 ( $r = 0.38$  and  $0.37$  for direct and reverse mutations respectively) compared to those of p53,  
341 though the bias is also low ( $r_{dir-rev} = -0.97$  and  $\langle\delta\rangle = -0.02$  kcal/mol). The poorer correlations  
342 for myoglobin are likely because the myoglobin data set consists of  $\Delta\Delta G$  measurements



343 obtained under various experimental conditions which ThermoNet does not explicitly account for.  
344 In fact, after excluding four data points (L29N, A130L, and two data points corresponding to  
345 A130K) with the biggest prediction error ( $> 3$  kcal / mol), the Pearson correlations increase to  
346 0.52 and 0.51 for direct and reverse mutations respectively. While the correlation between  
347 predicted and experimentally measured  $\Delta\Delta G$ s is not perfect, the predictions are generally  
348 conservative – no mutations with low measured  $\Delta\Delta G$  are predicted to have a high  $\Delta\Delta G$ . These  
349 results demonstrate the utility of ThermoNet as rapid unbiased predictor of  $\Delta\Delta G$  for mutations in  
350 clinically relevant proteins.

351 We also compared ThermoNet with four other biophysics-based methods: FoldX, Rosetta, SDM,  
352 and CUPSAT. We were not able to include all methods because both p53 and myoglobin are  
353 already in the S2648 set that was used for training most of the other machine learning-based  
354  $\Delta\Delta G$  predictors and they are also in the VariBench (36) data set used to derive parameters of  
355 the DDGun model (37). Our comparison indicates that both FoldX and Rosetta predictions have  
356 a better correlation than ThermoNet while also reasonably anti-symmetric (Supporting  
357 Information Table S7 and S8). However, as shown in Table 2 and demonstrated in previous  
358 studies (24, 25), both FoldX and Rosetta are likely to show bias toward predicting destabilization  
359 when tested on larger data sets.



360  
 361 Fig 4. Predicted ThermoNet  $\Delta\Delta G$  landscapes on p53 tumor suppressor protein and myoglobin.  
 362 (A) Performance of ThermoNet on predicting  $\Delta\Delta G$  for the direct mutations in p53 ( $r = 0.45$ ,  $\sigma =$   
 363  $2.01$  kcal/mol). (B) Performance of ThermoNet on predicting  $\Delta\Delta G$  for the reverse mutations in  
 364 p53 ( $r = 0.56$ ,  $\sigma = 1.92$  kcal/mol). (C) Direct versus reverse  $\Delta\Delta G$  values of all p53 mutations  
 365 predicted by ThermoNet ( $r = -0.93$  and  $\langle\delta\rangle = -0.04$  kcal/mol). (D) Performance of ThermoNet  
 366 on predicting  $\Delta\Delta G$  for the direct mutations in myoglobin ( $r = 0.38$ ,  $\sigma = 1.16$  kcal/mol). (E)  
 367 Performance of ThermoNet on predicting  $\Delta\Delta G$  for the reverse mutations in myoglobin ( $r = 0.37$ ,  
 368  $\sigma = 1.18$  kcal/mol). (F) Direct versus reverse  $\Delta\Delta G$  values of all myoglobin mutations predicted by  
 369 ThermoNet, with a Pearson correlation of  $r = -0.97$  and  $\langle\delta\rangle = -0.02$  kcal/mol. The dots are  
 370 colored in gradient from blue to red such that blue represents the most accurate prediction and  
 371 red indicates the least accurate prediction.

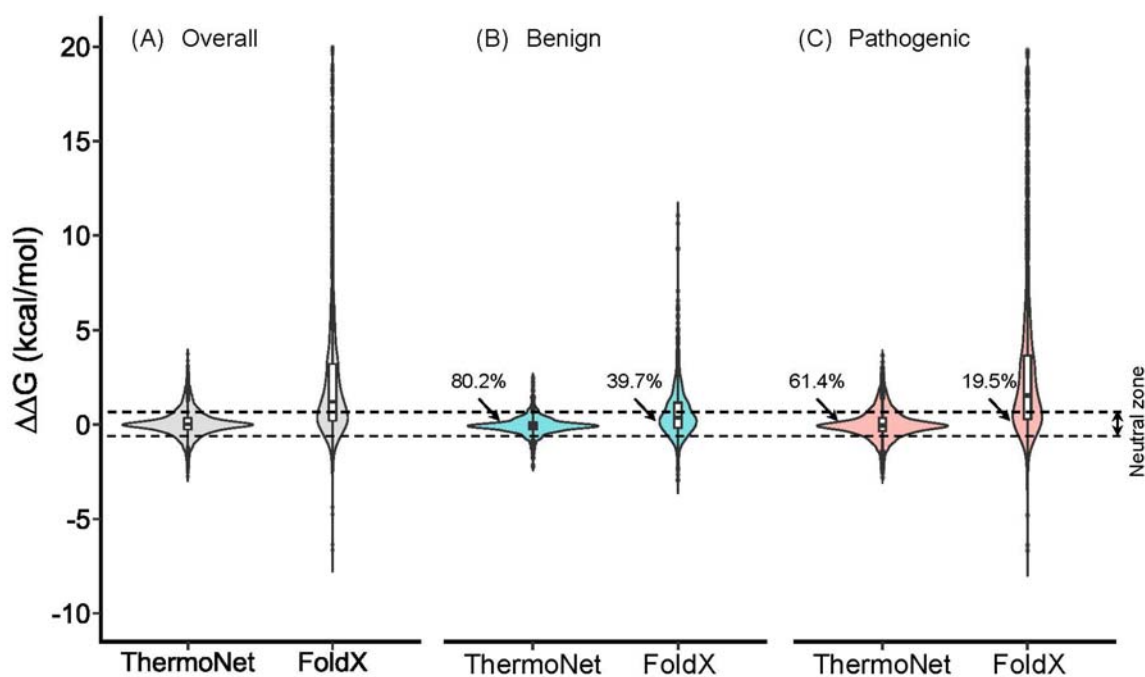
### 372 $\Delta\Delta G$ landscape of ClinVar missense variants

373 Previous work has shown that variant deleteriousness can only be partially attributed to  $\Delta\Delta G$  (43)  
 374 and that both stabilization and destabilization can cause disease (44). We sought to use  
 375 ThermoNet to obtain a less biased picture of the impact of benign and pathogenic variants on  
 376 proteins stability. We applied ThermoNet to predict the  $\Delta\Delta G$  distributions of pathogenic and

377 benign missense variants in ClinVar, a widely used resource of medically important variants (45).  
378 For comparison, we also applied FoldX, a popular and freely available  $\Delta\Delta G$  predictor (13, 22) to  
379 the ClinVar set. We first examined the overall predicted  $\Delta\Delta G$  distribution of ClinVar variants. The  
380  $\Delta\Delta G$ s of ClinVar variants predicted by ThermoNet range from -2.75 kcal/mol to +3.75 kcal/mol  
381 (Fig 5A). Experimentally measured  $\Delta\Delta G$ s generally fall within -5 kcal/mol to +5 kcal/mol (13, 46);  
382 thus, ThermoNet's predictions are consistent with the range of observed values. In contrast, the  
383  $\Delta\Delta G$ s of the same set of variants predicted by FoldX range from -6.64 kcal/mol to +57.2  
384 kcal/mol (Fig 5A), and 15.2% of  $\Delta\Delta G$ s predicted by FoldX are outside the expected range of -5  
385 kcal/mol to +5 kcal/mol.

386  $\Delta\Delta G$ s predicted by ThermoNet are also consistent with the expected  $\Delta\Delta G$  distributions  
387 according to a biophysical model of protein evolution (46). This model hypothesizes that fitness  
388 is a non-monotonic, concave function of protein stability, meaning that fitness decreases with  
389 increasing deviation from an optimal stability. The model also suggests that there is a "neutral  
390 zone" of 1.0 kcal/mol around the optimal stability in stability space, and mutations whose impact  
391 on stability fall within the neutral zone will have little effect on fitness (46). We thus reasoned  
392 that the  $\Delta\Delta G$ s of benign variants should fall within a narrow range from -0.5 kcal/mol to +0.5  
393 kcal/mol and that pathogenic variants should be equally likely to be destabilizing or stabilizing  
394 (46). To test this hypothesis, we examined the  $\Delta\Delta G$  distributions of pathogenic and benign  
395 variants separately. The  $\Delta\Delta G$ s of 80.2% of benign variants predicted by ThermoNet fall within  
396 the neutral zone, whereas FoldX only predicted 39.7% of benign variants to be in the neutral  
397 zone (Fig 5B). Further, the  $\Delta\Delta G$ s of pathogenic variants predicted by ThermoNet suggest  
398 pathogenic variants are nearly equally likely to be destabilizing (52.7%) as they are to be  
399 stabilizing (47.3%) (Fig 5C). In contrast, FoldX predicted that 83.2% of pathogenic variants are  
400 destabilizing. As already demonstrated in previous studies, the bias is likely because FoldX was  
401 parameterized on an experimental  $\Delta\Delta G$  data set dominated by destabilizing mutations (23-25).

402 ThermoNet's predictions are also consistent with the fact that variant pathogenicity can only be  
403 partially attributed to impacts on protein stability. ThermoNet predicts that 61.4% of pathogenic  
404 variants that have  $\Delta\Delta G$ s within the neutral zone and 19.8% of benign variants have  $\Delta\Delta G$ s  
405 outside the neutral zone (Fig 5B and 5C). This is expected based on previous biochemical  
406 characterizations of pathogenic mutations. For example, Bromberg *et al.* collected 66 mutations  
407 with experimentally measured  $\Delta\Delta G$  and functional annotations from the literature. The  $\Delta\Delta G$ s of  
408 this set of mutations range from -4.3 to 4.96 kcal/mol and the authors found that 31% of  
409 mutations affecting function had  $\Delta\Delta G$ s within the neutral zone while 19% functionally neutral  
410 mutations had  $\Delta\Delta G$ s outside the neutral zone (43).



411 Fig 5. Predicted  $\Delta\Delta G$  distributions of ClinVar missense variants.  
412 (A) The overall  $\Delta\Delta G$  distributions of ClinVar variants predicted by ThermoNet and FoldX.  
413 ThermoNet's predictions are consistent with the expected range based on experimentally  
414 determined  $\Delta\Delta G$  values (-5 kcal/mol to +5 kcal/mol). In contrast, more than 15% of  $\Delta\Delta G$ s  
415 predicted by FoldX are outside the expected range. (B) The  $\Delta\Delta G$  distributions for ClinVar benign  
416 variants predicted by ThermoNet and FoldX. (C) The  $\Delta\Delta G$  distributions of ClinVar pathogenic  
417 variants predicted by ThermoNet and FoldX. The  $\Delta\Delta G$ s of 80.2% of benign variants predicted by  
418

419 ThermoNet fall within the neutral zone (-0.5 to +0.5 kcal/mol, region between dashed lines), in  
420 which variants are not expected to influence fitness. FoldX only predicted 39.7% of benign  
421 variants to be in the neutral zone. Further, the  $\Delta\Delta G$ s of pathogenic variants predicted by  
422 ThermoNet suggest pathogenic variants are nearly equally likely to be stabilizing (47.3%) as  
423 destabilizing (52.7%). In contrast, FoldX predicted that 83.2% of pathogenic variants are  
424 destabilizing. Variants for which FoldX  $\Delta\Delta G$  is > 20 kcal/mol are omitted for clarity. Percentages  
425 represent the fractions of variants whose  $\Delta\Delta G$ s are predicted to be in the neutral zone.

## 426 **Discussion**

427 Accurate modeling of protein thermodynamic stability is a complex task due to the delicate  
428 balance between the different thermodynamic state functions that contribute to protein stability  
429 (1). The primary goal of this paper is to present a novel application of deep 3D convolutional  
430 neural networks to a fundamental challenge in structural bioinformatics: predicting changes in  
431 thermodynamic stability upon point mutation. We formulated the problem of  $\Delta\Delta G$  prediction from  
432 a computer vision perspective and took full advantage of the power of the constrained  
433 architecture of convolutional neural networks in detecting spatially proximate features (29). We  
434 developed ThermoNet, a method based on deep 3D convolutional neural networks, to predict  
435  $\Delta\Delta G$  upon point mutation. We showed that  $\Delta\Delta G$  can be predicted from protein structure with  
436 reasonable accuracy using deep 3D convolutional neural networks without manual feature  
437 engineering. While ThermoNet achieved comparable performance to previous methods on  
438 direct mutations, it performed better on reverse mutations than most methods by a large margin,  
439 and remarkably, reduced the magnitude of prediction bias.

440 In addition to introducing ThermoNet, we also address two methodological challenges in the  
441 development and evaluation of computational methods for  $\Delta\Delta G$  prediction: lack of anti-  
442 symmetry and data leaks due to homology between proteins in the training and evaluation sets.  
443 Previously, it was shown that the lack of anti-symmetry in  $\Delta\Delta G$  prediction can be effectively  
444 addressed either by using input features that are anti-symmetric by construction (19, 37, 38, 47,  
445 48) or by training the predictor using both direct and reverse mutations (19, 24). In addition,  
446 when the statistical model is parametric, one may identify the terms that are responsible for

447 breaking the symmetry and make correction accordingly (48). However, when the predictor is  
448 nonparametric, meaning that the terms of the statistical learning model are not established a  
449 priori, the only way is to train the model with data set balanced with direct and reverse mutations  
450 so that it learns the anti-symmetry (24). For structure-based  $\Delta\Delta G$  predictors, this approach  
451 requires knowing the 3D structures of both the wild type and mutant protein. While mutant  
452 structures determined via experimental techniques are scarce, in this work, we demonstrated  
453 that mutant protein structures obtained through molecular modeling-based data augmentation  
454 can also be effectively used as substitutes for experimental structures to remedy the lack of  
455 anti-symmetry in  $\Delta\Delta G$  prediction.

456 Recently, the potential for data leak between training and testing due to the inclusion of  
457 mutations from the same proteins has been appreciated (16, 19); however, the effects of  
458 including mutations from homologous proteins in training and validation sets is less appreciated  
459 and understood. Our results suggest that that such homology can influence performance  
460 estimates. The ThermoNet\* model, which was trained before homology reduction, achieved  
461 stronger performance than the ThermoNet model trained after homology reduction (Fig 2B). In  
462 real-world applications, there will be homology between proteins used to train prediction models  
463 and the proteins to which they are applied. However, given the relatively small number of  
464 proteins included in commonly used training sets and the fact that they are not representative of  
465 the full diversity of protein folds and functions, we believe that the inclusion of mutations from  
466 proteins with shared evolutionary histories is likely to bias performance estimates. In the future,  
467 it will be valuable to explore this issue further and construct training sets that reflect the  
468 evolutionary relationships expected in various applications.

469 ThermoNet treats protein structures as if they were 3D images, and it takes as input a 3D grid of  
470 voxels parameterized with seven biophysical property channels. As such, this approach  
471 bypasses the tedious processes of manual feature engineering and feature selection that, if not

472 done correctly, can often lead to over-optimistic estimation of model performance. The locally  
473 constrained deep convolutional architecture likely allows the system to model the complex, non-  
474 linear nature of molecular interactions. Recently, a spherical convolutional architecture in which  
475 concentric voxel grids parameterized by atom masses and charges were used as input to  
476 predict the  $\Delta\Delta G$ s of direct mutations with good accuracy (49). Thus, together with the current  
477 work, these results demonstrate the potential of deep convolutional neural networks for  
478 predicting biophysically meaningful information from protein structures and holds promise for  
479 protein engineering.

480 The fact that our approach relies on the availability of experimental structures or homology  
481 models and the 3D nature of the convolutional neural network create two limitations. First, while  
482 protein structures are being determined at an unprecedented pace, the fraction of the human  
483 proteome with available experimental structure is estimated to be around 20% (50). Even when  
484 all the proteins whose structures can be modeled reliably are considered, only ~70% of the  
485 human proteome will have structural coverage (50). As the structures of many proteins can only  
486 be partially modeled, the space of the human proteome that one can apply ThermoNet to will be  
487 less than 70%. Second, compared to the 2D version with the same architecture, ThermoNet has  
488 four times more parameters: three convolutional layers of 16, 24, and 32 neurons respectively,  
489 and one dense layer with 24 neurons that takes an input tensor of the shape [16, 16, 16, 14] has  
490 133,273 parameters. Training deep 3D CNNs is very demanding, requiring more GPU memory  
491 and more training data to avoid overfitting. While we demonstrated the potential of deep 3D  
492 CNNs in modeling  $\Delta\Delta G$  of proteins, the relatively little training data available raises the question  
493 of whether deep 3D CNNs can model  $\Delta\Delta G$ s and related thermodynamic properties at  
494 experimental accuracy. Nonetheless, the increasing adoption of deep mutational scanning  
495 techniques for systematic study of the molecular effects of mutations is generating an  
496 unprecedented amount of data. Furthermore, given rapid increases in GPU power and the

497 number of structures of proteins and their complexes determined, we expect that deep 3D  
498 CNNs will be successfully applied to provide solutions to many biophysical problems such as  
499 modeling the impact of mutation on protein-protein, protein-DNA as well as protein-RNA  
500 interactions.

## 501 **Methods**

### 502 ***Protein thermodynamic stability***

503 The thermodynamic stability  $\Delta G$  of the folded form of a two-state protein, which is its Gibbs free  
504 energy of folding, is defined in relation to the concentration of folded [*folded*] and the  
505 concentration of unfolded [*unfolded*] forms:

$$\Delta G = -RT \ln \frac{[folded]}{[unfolded]}$$

506 where T is the temperature, and R is the gas constant. More stable proteins, meaning that a  
507 higher fraction of the protein is in the folded form, have more negative values of  $\Delta G$ . The impact  
508 of mutations on protein stability,  $\Delta\Delta G$ , is defined in terms of the change in  $\Delta G$  between the wild-  
509 type and mutant proteins:

$$\Delta\Delta G = \Delta G_{mutant} - \Delta G_{wild-type} = -RT \ln \frac{[folded]_{mutant}/[unfolded]_{mutant}}{[folded]_{wild-type}/[unfolded]_{wild-type}}$$

510 such that a destabilizing mutation has a positive  $\Delta\Delta G$ , whereas a stabilizing mutation has a  
511 negative  $\Delta\Delta G$ . The values of  $\Delta\Delta G$  resulting from single-point mutations usually range from -5 to  
512 5 kcal/mol (13). A mutation that destabilizes a typical protein ( $\Delta G = -5$  kcal/mol) by 1 kcal/mol  
513 will reduce the equilibrium constant for the folding reaction of this protein by a factor of 5.1 at  
514 physiological temperature (310 K).



## 515 ***Thermodynamics of direct and reverse mutations***

516 Consider a pair of proteins whose sequences differ only at a single position where the amino  
517 acid is  $X$  in one protein and  $Y$  in the other. Let the Gibbs free energies of folding of this pair of  
518 proteins be  $\Delta G_X$  and  $\Delta G_Y$  respectively. For such a pair of proteins, one can think of the protein  $Y$   
519 as being generated by a “direct” mutation at the sequence location from amino acid  $X$  to  $Y$  and  
520 the change in the Gibbs free energy of folding caused by  $X$  to  $Y$  mutation is  $\Delta\Delta G_{X\rightarrow Y} = \Delta G_Y -$   
521  $\Delta G_X$ . One may also think of the protein  $X$  as being generated by a “reverse” mutation from  
522 amino acid  $Y$  to  $X$  and the change in the Gibbs free energy of folding caused by this reverse  
523 mutation is  $\Delta\Delta G_{Y\rightarrow X} = \Delta G_X - \Delta G_Y = -\Delta\Delta G_{X\rightarrow Y}$ . A well-performing, “self-consistent” method for  
524 predicting  $\Delta\Delta G$ s would not only give accurate  $\Delta\Delta G$  predictions for the direct mutations, but also  
525 for the reverse mutations. The self-consistency requirement has been largely ignored by  
526 previously developed  $\Delta\Delta G$  predictors (23-25).

## 527 ***Data sets and symmetry-based data augmentation***

528 The data set used to train ThermoNet was derived from the Q3421 data set compiled in a  
529 previous study (15). The Q3421 data set contains 3,421 distinct single-point mutations in 150  
530 proteins collected from the ProTherm database (51). The impact of these mutations on the  
531 stability of the protein structure have been measured experimentally and expressed  
532 quantitatively as  $\Delta\Delta G$  values. We first excluded those mutations from the Q3421 data set that  
533 were also in the  $S^{\text{sym}}$  test set (see following). To reduce the sequence similarity between  
534 proteins in the training set of ThermoNet and the proteins it was tested on, we also removed all  
535 proteins that are likely homologous (BLAST evalue  $< 0.001$ ) to p53, myoglobin, and proteins in  
536  $S^{\text{sym}}$  from Q3421. Estimation of homology was accomplished by running the blastp program (52)  
537 using protein sequences in the  $S^{\text{sym}}$  data set and the sequences of p53 and myoglobin as  
538 queries against protein sequences in Q3421. Our rigorous pruning of the Q3421 data set  
539 resulted in a final data set consisting of 1,744 distinct mutations in 127 proteins. This data set

540 was augmented by creating a reverse mutation data point for each of the 1,744 direct mutations,  
541 thus giving to a total of 3,488 data points for the training of ThermoNet. The data set was  
542 randomly divided into ten equally sized, mutually exclusive subsets each consisting of 10% of  
543 the direct mutations and the corresponding 10% of reverse mutations. In the training of each  
544 component model of ThermoNet, nine subsets were combined to form a training set and the  
545 remaining one subset was used as a validation set. The data set used to test ThermoNet and to  
546 compare it with fifteen previously developed  $\Delta\Delta G$  predictors was a common, balanced data set  
547 called  $S^{\text{sym}}$  consisting of 342 pairs of proteins with known crystal structures (24). The members  
548 forming each pair differ at only a single position in the protein sequence. The  $\Delta\Delta G$  values of the  
549 342 direct mutations have been experimentally measured and the  $\Delta\Delta G$  values of the  
550 corresponding 342 reverse mutations were assigned using anti-symmetry.

### 551 ***Modeling mutant structures***

552 We treat each mutation as a pair of proteins whose sequences differ only at a single sequence  
553 position. For each pair of proteins, we designate the one whose structure has been  
554 experimentally resolved as protein X and the other as protein Y. Structures of the X proteins  
555 were collected from the Protein Data Bank (53) and were relaxed in the Rosetta all-atom energy  
556 function ref2015 (54) using the Rosetta *FastRelax* protocol (42). To prevent large-scale  
557 conformational shift from the input PDB structure, atoms were constrained to their starting  
558 locations with a harmonic penalty potential during relaxation. The same Rosetta *FastRelax*  
559 protocol was also employed to create structural models for each of the Y proteins from the  
560 corresponding relaxed structure of the X protein by supplying a Rosetta *resfile* specifying the  
561 mutation  $X \rightarrow Y$  to make. The structures of both proteins of each protein pair in the test set were  
562 collected from the Protein Data Bank (53) and were also relaxed using the same Rosetta  
563 *FastRelax* protocol.

## 564 ***Voxelization of the neighborhood of mutation site***

565 We treated each protein structure as a collection of volume elements (voxels) in 3D space. Just  
566 as a pixel element in an image has color channels, we parameterized a voxel in a protein  
567 structure by a set of  $k$  chemical property channels:  $[v_1, v_2, \dots, v_k]$  where the value  $v_i$  of each  
568 property channel indicates the level of saturation of property  $i$  at this voxel (Fig 1A and 1B). For  
569 each mutation (a pair of proteins), we superimposed the mutant structure onto the wild-type  
570 structure such that the root-mean-squared distance between them is minimized and collected a  
571 grid of  $16 \times 16 \times 16$  voxels from both structures. We parameterized each voxel with seven  
572 property channels each of a distinct chemical nature according to AutoDock4 atom types (55) as  
573 in the work of Jimenez et al. (28) (Table 1). This resulted in a tensor of the shape  $[16, 16, 16, 7]$   
574 for a single structure and a tensor of the shape  $[16, 16, 16, 14]$  when the two tensors from both  
575 structures are concatenated to represent the mutation. The grid was centered at the  $C_\beta$  atom of  
576 the mutation site amino acid (or  $C_\alpha$  atom if it's a glycine) where each voxel is a unit cube whose  
577 sides are 1 Å long. The level of saturation  $f(d)$  of each property channel at each voxel is  
578 determined by the van der Waals radius  $r_{vdw}$  of the atom designated to have that property and  
579 its distance  $d$  to the center of the voxel through the following formula:

$$f(d) = 1 - \exp \left[ - \left( \frac{r_{vdw}}{d} \right)^{12} \right]$$

580 The computation of tensors from protein structures was performed using routines implemented  
581 in the HTMD Python library (version 1.17) for molecular simulations (56) and a Python program  
582 for creating a data set from a list of mutations is provided in Supplementary Code.

## 583 ***Model architecture***

584 Convolutional neural networks are a type of deep-learning model commonly used in computer  
585 vision applications. They have recently proven to perform well in residue contact prediction (57-  
586 60) and protein tertiary structure prediction (61-64). We selected CNNs for protein structure-

587 based  $\Delta\Delta G$  prediction because we formulated this problem as a computer vision problem by  
588 treating protein structures as if they were 3D images. Each of the component models of  
589 ThermoNet features a sequential organization of three 3D convolutional layers (Conv3D), one  
590 3D max pooling layer (MaxPool3D), followed by one fully connected layer (Dense) (Fig S1).  
591 Convolutions operate over 4D tensors, called feature maps, with three spatial axes (length,  
592 width, and height) as well as a depth axis. The convolution operation extracts 3D patches of  
593 shape [3, 3, 3] with stride 1 from its input feature map and applies the same transformation to all  
594 patches, producing an output feature map. This output feature map is still a 4D tensor: it has a  
595 length, height, and a width whose values are determined by the shape of convolution patches  
596 and size of the stride, but its depth, which is also called number of filters, is a hyperparameter of  
597 the layer (see the section on hyperparameter search below). The number of filters in the three  
598 Conv3D layers were 16, 24, and 32 respectively. All convolution operation outputs from each  
599 Conv3D layer are transformed by the rectified linear activation function (ReLU). The  
600 transformed outputs from the last Conv3D layer are pooled by taking the maximum activation of  
601 each  $2 \times 2 \times 2$  grid. The max-pooled activations are then flattened into a 1D vector of features  
602 which are fully connected with a dense layer of 24 ReLU units. This model architecture was  
603 implemented using Keras (65) with TensorFlow (66) as the backend.

#### 604 ***Training ThermoNet***

605 ThermoNet is an ensemble of ten deep 3D convolutional neural networks, each trained to  
606 perform the best on a validation set. Each of the component models was trained on nine  
607 subsets (collectively known as the training set), and its generalization performance was  
608 monitored on the remaining one subset (validation set). This process was iterated ten times  
609 each using a different one of the ten subsets as the validation set and the remaining nine  
610 subsets as the training set. We employed this procedure to obtain an ensemble of ten models  
611 as model ensembling has been suggested to produce better predictions (67). Evaluation of this

612 ensemble was performed on a separate test set (see below). All component models of  
613 ThermoNet were trained using the Adam optimizer (68) for 200 epochs with default  
614 hyperparameters (maximum learning rate = 0.001,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ). Kernel weights of the  
615 model were initialized using the Glorot uniform initializer and updated after processing each  
616 batch of eight training examples. The mean squared error (MSE) of the predicted  $\Delta\Delta G$  values  
617 from the experimental measurements was used as the loss function during training. When  
618 training a deep neural network, one often cannot predict how many epochs will be needed to get  
619 to an optimal validation loss. We monitored the MSE of the predictions on a separate validation  
620 set consisting of 10% variants randomly selected from the training set during training. Training  
621 was stopped when the MSE on the validation set stopped decreasing for ten consecutive  
622 epochs. To regularize the model, the dense layer was placed between two dropout layers with a  
623 dropout rate of 0.5 in each layer (Fig S1). Each of the final component models was the one that  
624 produced the lowest MSE on the validation set. The final predicted  $\Delta\Delta G$  value is the average of  
625 predictions from the ten models.

### 626 ***Hyperparameter search***

627 The design of deep CNNs entails many architectural choices to account for number of hidden  
628 convolutional layers and fully connected layers, number of filters, filter size, strides, padding,  
629 dropout rate among many other hyperparameters. We initially created a voxel grid with size  
630  $16 \times 16 \times 16$  at a resolution of 1 Å for each chemical property channel following the procedure  
631 described in (27) to cross-validate our network architectures. Considering the limited training  
632 data set available in our study, we tried some smaller architectures with cross-validation to  
633 decide the optimal one, rather than simply adapting the widely used, much larger, network  
634 architectures in computer vision applications. We restricted our deep CNNs to have three  
635 convolutional layers and one fully connected while considering several hidden layer sizes. Our  
636 results from five-fold cross-validation suggest that the architecture of the  $16 \times 24 \times 32$

637 convolutional configuration combined with a fully connected layer of size 24 achieved the best  
638 performance (Fig S1). An additional consideration for our deep 3D CNN architecture is the  
639 dimension of the local box. The size of the local box specifies the structural information  
640 accessible by the network and therefore is a hyperparameter of our method. We cross-validated  
641 the voxelization scheme with grid of sizes  $8 \times 8 \times 8$ ,  $12 \times 12 \times 12$ ,  $16 \times 16 \times 16$ , and  $20 \times 20 \times$   
642  $20$  at a resolution of  $1 \text{ \AA}$ . All voxelization schemes draw a cubic box around the mutation site  
643 with lateral length of  $l \text{ \AA}$  where  $l$  equals 8, 12, 16, or  $20 \text{ \AA}$ . Our results from five-fold cross-  
644 validation indicate that the voxel grid with size  $16 \times 16 \times 16$  gives the best prediction  
645 performance (Fig S1).

#### 646 **Performance evaluation**

647 The following measures were adopted to evaluate the performance of ThermoNet and to  
648 facilitate comparison with previously developed methods. The primary measures for evaluating  
649 prediction accuracy were the Pearson correlation coefficient ( $r$ ) between experimental and  
650 predicted  $\Delta\Delta G$ s and the root-mean-squared error ( $\sigma$ ) of predictions. For a set of  $n$  data points  
651  $(x_i, y_i)$ , the formula for calculating  $r$  and  $\sigma$  are defined as follows:

$$r = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - x_i)^2}$$

652 where the  $(x_i, y_i)$  tuple denotes the experimental and predicted  $\Delta\Delta G$  values of mutation  $i$ ,  
653 respectively, and  $n$  denotes the number of mutations in the data set. The measures for  
654 evaluating prediction bias were the Pearson correlation coefficient between the predictions for  
655 direct mutations and those for reverse mutations and the parameter  $\delta$  which is defined as:

656  $\delta = \Delta\Delta G_{inv} + \Delta\Delta G_{dir}$  and was previously used to quantify prediction bias (23). An unbiased  
657 predictor should have  $\delta = 0$  for every mutation. The average of  $\delta$ ,  $\langle\delta\rangle$ , taken over all mutations  
658 in the  $S^{\text{sym}}$  data set was used in two previously studies to quantify prediction bias (24, 37). While  
659 we report  $\langle\delta\rangle$  in this work, we note that  $\langle\delta\rangle$  is flawed because biases toward opposite directions  
660 will be washed out when summed. To give a more transparent presentation of prediction bias,  
661 we also plot the distribution of  $\delta$  and report the average of absolute bias, i.e.  $\langle|\delta|\rangle$ .

### 662 ***ClinVar variants***

663 We retrieved the ClinVar database (45) in VCF format on August 15, 2019 and ran the VCF file  
664 through the Variant Effect Predictor (version 97) (69) to annotate the consequences of all  
665 ClinVar variants. We created a set of missense variants that can be mapped to protein  
666 structures to demonstrate the applicability of ThermoNet to clinically relevant variants. Our  
667 evaluation set consists of solely ClinVar missense variants that are labeled as "pathogenic" or  
668 "likely pathogenic" for true positive (pathogenic) variants and "benign" or "likely benign" for true  
669 negative (benign) variants. All variants are required to have a review status of at least one star  
670 and no conflicting interpretation. Any ClinVar variant designated as "no assertion criteria  
671 provided", "no assertion provided", "no interpretation for the single variant", or not covered by  
672 either a structure or homology model was excluded from the evaluation set. Due to the  
673 dependency of ThermoNet on 3D structures, we also require variants in the evaluation set to be  
674 mappable to available protein structures. The residue-level mapping of ClinVar variants onto  
675 protein structures was based on the SIFTS resource that provides residue-level mapping  
676 between UniProt and Protein Data Bank (PDB) entries (70). Collectively, these restrictions  
677 resulted in 3,510 pathogenic variants and 950 benign variants that can be mapped to  
678 experimental structures deposited in the PDB (53). The mapped variants along with PDB IDs  
679 can be found in our GitHub repository at <https://github.com/gersteinlab/ThermoNet>.

680 **Author contribution**

681 B.L. and M.B.G. conceived the study. B.L. implemented ThermoNet. B.L. and Y.T.Y performed  
682 the analysis. B.L. performed the work on predicting the  $\Delta\Delta G$  distributions of ClinVar variants  
683 using ThermoNet at Vanderbilt University with the help from J.A.C. M.B.G. and J.A.C.  
684 supervised the project. B.L. wrote the manuscript with help from all authors.

685 **Code and data availability**

686 ThermoNet source code and raw data supporting the analysis of this work is available at  
687 <https://github.com/gersteinlab/ThermoNet>.

688 **Acknowledgement**

689 This work was supported by NSF award DBI1660648 (M.B.G.), NIH awards R35 GM127087  
690 (J.A.C.) and R01 GM126249 (J.A.C.), and an American Heart Association Postdoctoral  
691 Fellowship 20POST35220002 (B.L.).The authors would also like to thank the Center for  
692 Research Computing at Yale University and the Advanced Computing Center for Research and  
693 Education at Vanderbilt University for supporting high-performance computing.

694 **References**

- 695 1. Li B, Fooksa M, Heinze S, Meiler J. Finding the needle in the haystack: towards solving  
696 the protein-folding problem computationally. *Crit Rev Biochem Mol Biol*. 2018;53(1):1-28.
- 697 2. Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, et al. A global  
698 reference for human genetic variation. *Nature*. 2015;526(7571):68-74.
- 699 3. Wang Z, Moulton J. SNPs, protein structure, and disease. *Hum Mutat*. 2001;17(4):263-70.
- 700 4. Yue P, Li Z, Moulton J. Loss of protein structure stability as a major causative factor in  
701 monogenic disease. *J Mol Biol*. 2005;353(2):459-73.
- 702 5. Stein A, Fowler DM, Hartmann-Petersen R, Lindorff-Larsen K. Biophysical and  
703 Mechanistic Models for Disease-Causing Protein Variants. *Trends Biochem Sci*.  
704 2019;44(7):575-88.
- 705 6. Huang PS, Boyken SE, Baker D. The coming of age of de novo protein design. *Nature*.  
706 2016;537(7620):320-7.
- 707 7. Gapsys V, Michielssens S, Seeliger D, de Groot BL. Accurate and Rigorous Prediction  
708 of the Changes in Protein Free Energies in a Large-Scale Mutation Scan. *Angew Chem Int Edit*.  
709 2016;55(26):7364-8.
- 710 8. Kellogg EH, Leaver-Fay A, Baker D. Role of conformational sampling in computing  
711 mutation-induced changes in protein structure and stability. *Proteins-Structure Function and*  
712 *Bioinformatics*. 2011;79(3):830-8.



- 713 9. Bender BJ, Cisneros A, Duran AM, Finn JA, Fu D, Lokits AD, et al. Protocols for  
714 Molecular Modeling with Rosetta3 and RosettaScripts. *Biochemistry*. 2016;55(34):4748-63.
- 715 10. Yin SY, Ding F, Dokholyan NV. Eris: an automated estimator of protein stability. *Nature*  
716 *Methods*. 2007;4(6):466-7.
- 717 11. Worth CL, Preissner R, Blundell TL. SDM-a server for predicting effects of mutations on  
718 protein stability and malfunction. *Nucleic Acids Research*. 2011;39:W215-W22.
- 719 12. Dehouck Y, Kwasigroch JM, Gilis D, Rooman M. PoPMuSiC 2.1: a web server for the  
720 estimation of protein stability changes upon mutation and sequence optimality. *BMC*  
721 *bioinformatics*. 2011;12.
- 722 13. Guerois R, Nielsen JE, Serrano L. Predicting changes in the stability of proteins and  
723 protein complexes: A study of more than 1000 mutations. *Journal of Molecular Biology*.  
724 2002;320(2):369-87.
- 725 14. Parthiban V, Gromiha MM, Schomburg D. CUPSAT: prediction of protein stability upon  
726 point mutations. *Nucleic Acids Research*. 2006;34:W239-W42.
- 727 15. Quan LJ, Lv Q, Zhang Y. STRUM: structure-based prediction of protein stability changes  
728 upon single-point mutation. *Bioinformatics*. 2016;32(19):2936-46.
- 729 16. Pires DEV, Ascher DB, Blundell TL. mCSM: predicting the effects of mutations in  
730 proteins using graph-based signatures. *Bioinformatics*. 2014;30(3):335-42.
- 731 17. Capriotti E, Fariselli P, Casadio R. I-Mutant2.0: predicting stability changes upon  
732 mutation from the protein sequence or structure. *Nucleic Acids Research*. 2005;33:W306-W10.
- 733 18. Masso M, Vaisman II. Accurate prediction of stability changes in protein mutants by  
734 combining machine learning with structure based computational mutagenesis. *Bioinformatics*.  
735 2008;24(18):2002-9.
- 736 19. Fariselli P, Martelli PL, Savojardo C, Casadio R. INPS: predicting the impact of non-  
737 synonymous variations on protein stability from sequence. *Bioinformatics*. 2015;31(17):2816-21.
- 738 20. Cao H, Wang J, He L, Qi Y, Zhang JZ. DeepDDG: Predicting the Stability Change of  
739 Protein Point Mutations Using Neural Networks. *Journal of Chemical Information and Modeling*.  
740 2019.
- 741 21. Roushar FJ, Gruenhagen TC, Penn WD, Li B, Meiler J, Jastrzebska B, et al.  
742 Contribution of Cotranslational Folding Defects to Membrane Protein Homeostasis. *J Am Chem*  
743 *Soc*. 2019;141(1):204-15.
- 744 22. Buss O, Rudat J, Ochsenreither K. FoldX as Protein Engineering Tool: Better Than  
745 Random Based Approaches? *Comput Struct Biotechnol J*. 2018;16:25-33.
- 746 23. Thiltgen G, Goldstein RA. Assessing Predictors of Changes in Protein Stability upon  
747 Mutation Using Self-Consistency. *Plos One*. 2012;7(10).
- 748 24. Pucci F, Bernaerts KV, Kwasigroch JM, Rooman M. Quantification of biases in  
749 predictions of protein stability changes upon mutations. *Bioinformatics*. 2018;34(21):3659-65.
- 750 25. Usmanova DR, Bogatyreva NS, Bernad JA, Eremina AA, Gorshkova AA, Kanevskiy GM,  
751 et al. Self-consistency test reveals systematic bias in programs for prediction change of stability  
752 upon mutation. *Bioinformatics*. 2018;34(21):3653-8.
- 753 26. Fang J. A critical review of five machine learning-based algorithms for predicting protein  
754 stability changes upon mutation. *Brief Bioinform*. 2019.
- 755 27. Jimenez J, Doerr S, Martinez-Rosell G, Rose AS, De Fabritiis G. DeepSite: protein-  
756 binding site predictor using 3D-convolutional neural networks. *Bioinformatics*.  
757 2017;33(19):3036-42.
- 758 28. Jimenez J, Skalic M, Martinez-Rosell G, De Fabritiis G. KDEEP: Protein-Ligand Absolute  
759 Binding Affinity Prediction via 3D-Convolutional Neural Networks. *J Chem Inf Model*.  
760 2018;58(2):287-96.
- 761 29. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature*. 2015;521(7553):436-44.
- 762 30. Torng W, Altman RB. 3D deep convolutional neural networks for amino acid  
763 environment similarity analysis. *BMC bioinformatics*. 2017;18(1):302.

- 764 31. Torng W, Altman RB. High precision protein functional site detection using 3D  
765 convolutional neural networks. *Bioinformatics*. 2019;35(9):1503-12.
- 766 32. Wallach I, Dzamba M, Heifets A. AtomNet: A Deep Convolutional Neural Network for  
767 Bioactivity Prediction in Structure-based Drug Discovery. *arXiv*. 2015.
- 768 33. Leaver-Fay A, Tyka M, Lewis SM, Lange OF, Thompson J, Jacak R, et al. ROSETTA3:  
769 an object-oriented software suite for the simulation and design of macromolecules. *Methods*  
770 *Enzymol*. 2011;487:545-74.
- 771 34. Dehouck Y, Grosfils A, Folch B, Gilis D, Bogaerts P, Rooman M. Fast and accurate  
772 predictions of protein stability changes upon mutations using statistical potentials and neural  
773 networks: PoPMuSiC-2.0. *Bioinformatics*. 2009;25(19):2537-43.
- 774 35. Wang G, Dunbrack RL, Jr. PISCES: a protein sequence culling server. *Bioinformatics*.  
775 2003;19(12):1589-91.
- 776 36. Yang Y, Urolagin S, Niroula A, Ding X, Shen B, Vihinen M. PON-tstab: Protein Variant  
777 Stability Predictor. Importance of Training Data Quality. *Int J Mol Sci*. 2018;19(4).
- 778 37. Montanucci L, Capriotti E, Frank Y, Ben-Tal N, Fariselli P. DDGun: an untrained method  
779 for the prediction of protein stability changes upon single and multiple point variations. *BMC*  
780 *bioinformatics*. 2019;20(Suppl 14):335.
- 781 38. Montanucci L, Savojardo C, Martelli PL, Casadio R, Fariselli P. On the biases in  
782 predictions of protein stability changes upon variations: the INPS test case. *Bioinformatics*.  
783 2019;35(14):2525-7.
- 784 39. Pucci F, Bourgeas R, Rooman M. High-quality Thermodynamic Data on the Stability  
785 Changes of Proteins Upon Single-site Mutations. *Journal of Physical and Chemical Reference*  
786 *Data*. 2016;45(2).
- 787 40. Ordway GA, Garry DJ. Myoglobin: an essential hemoprotein in striated muscle. *J Exp*  
788 *Biol*. 2004;207(Pt 20):3441-6.
- 789 41. Kepp KP. Towards a "Golden Standard" for computing globin stability: Stability and  
790 structure sensitivity of myoglobin mutants. *Biochimica et biophysica acta*. 2015;1854(10 Pt  
791 A):1239-48.
- 792 42. Tyka MD, Keedy DA, Andre I, Dimaio F, Song Y, Richardson DC, et al. Alternate states  
793 of proteins revealed by detailed energy landscape mapping. *J Mol Biol*. 2011;405(2):607-18.
- 794 43. Bromberg Y, Rost B. Correlating protein function and stability through the analysis of  
795 single amino acid substitutions. *BMC bioinformatics*. 2009;10 Suppl 8:S8.
- 796 44. Ancien F, Pucci F, Godfroid M, Rooman M. Prediction and interpretation of deleterious  
797 coding variants in terms of protein structural stability. *Sci Rep*. 2018;8(1):4480.
- 798 45. Landrum MJ, Lee JM, Benson M, Brown GR, Chao C, Chitipiralla S, et al. ClinVar:  
799 improving access to variant interpretations and supporting evidence. *Nucleic Acids Res*.  
800 2018;46(D1):D1062-D7.
- 801 46. DePristo MA, Weinreich DM, Hartl DL. Missense meanderings in sequence space: a  
802 biophysical view of protein evolution. *Nat Rev Genet*. 2005;6(9):678-87.
- 803 47. Savojardo C, Martelli PL, Casadio R, Fariselli P. On the critical review of five machine  
804 learning-based algorithms for predicting protein stability changes upon mutation. *Brief Bioinform*.  
805 2019.
- 806 48. Pucci F, Bernaerts K, Teheuxa F, Gilisa D, Roomana M. Symmetry Principles in  
807 Optimization Problems: an application to Protein Stability Prediction. *IFAC-PapersOnLine*.  
808 2015;48(1):458-63.
- 809 49. Boomsma W, Frellsen J, editors. Spherical convolutions and their application in  
810 molecular modelling. *Advances in Neural Information Processing Systems*; 2017.
- 811 50. Somody JC, MacKinnon SS, Windemuth A. Structural coverage of the proteome for  
812 pharmaceutical applications. *Drug Discov Today*. 2017;22(12):1792-9.

- 813 51. Kumar MDS, Bava KA, Gromiha MM, Prabakaran P, Kitajima K, Uedaira H, et al.  
814 ProTherm and ProNIT: thermodynamic databases for proteins and protein-nucleic acid  
815 interactions. *Nucleic Acids Research*. 2006;34:D204-D6.
- 816 52. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool.  
817 *J Mol Biol*. 1990;215(3):403-10.
- 818 53. Rose PW, Prlic A, Altunkaya A, Bi CX, Bradley AR, Christie CH, et al. The RCSB protein  
819 data bank: integrative view of protein, gene and 3D structural information. *Nucleic Acids*  
820 *Research*. 2017;45(D1):D271-D81.
- 821 54. Alford RF, Leaver-Fay A, Jeliazkov JR, O'Meara MJ, DiMaio FP, Park H, et al. The  
822 Rosetta All-Atom Energy Function for Macromolecular Modeling and Design. *J Chem Theory*  
823 *Comput*. 2017;13(6):3031-48.
- 824 55. Morris GM, Huey R, Lindstrom W, Sanner MF, Belew RK, Goodsell DS, et al. AutoDock4  
825 and AutoDockTools4: Automated Docking with Selective Receptor Flexibility. *Journal of*  
826 *Computational Chemistry*. 2009;30(16):2785-91.
- 827 56. Doerr S, Harvey MJ, Noe F, De Fabritiis G. HTMD: High-Throughput Molecular  
828 Dynamics for Molecular Discovery. *Journal of Chemical Theory and Computation*.  
829 2016;12(4):1845-52.
- 830 57. Kandathil SM, Greener JG, Jones DT. Prediction of interresidue contacts with  
831 DeepMetaPSICOV in CASP13. *Proteins*. 2019.
- 832 58. Schaarschmidt J, Monastyrskyy B, Kryshchuk A, Bonvin A. Assessment of contact  
833 predictions in CASP12: Co-evolution and deep learning coming of age. *Proteins*. 2018;86 Suppl  
834 1:51-66.
- 835 59. Wang S, Sun SQ, Xu JB. Analysis of deep learning methods for blind protein contact  
836 prediction in CASP12. *Proteins-Structure Function and Bioinformatics*. 2018;86:67-77.
- 837 60. Shrestha R, Fajardo E, Gil N, Fidelis K, Kryshchuk A, Monastyrskyy B, et al.  
838 Assessing the accuracy of contact predictions in CASP13. *Proteins*. 2019.
- 839 61. Xu J. Distance-based protein folding powered by deep learning. *Proc Natl Acad Sci U S*  
840 *A*. 2019.
- 841 62. Xu J, Wang S. Analysis of distance-based protein structure prediction by deep learning  
842 in CASP13. *Proteins*. 2019.
- 843 63. Kandathil SM, Greener JG, Jones DT. Recent developments in deep learning applied to  
844 protein structure prediction. *Proteins*. 2019.
- 845 64. Greener JG, Kandathil SM, Jones DT. Deep learning extends de novo protein modelling  
846 coverage of genomes using iteratively predicted structural constraints. *Nat Commun*.  
847 2019;10(1):3977.
- 848 65. Chollet F. keras. [url{https://github.com/fchollet/keras}](https://github.com/fchollet/keras); 2015.
- 849 66. Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, et al. TensorFlow: a system for  
850 large-scale machine learning. *Proceedings of the 12th USENIX conference on Operating*  
851 *Systems Design and Implementation*; Savannah, GA, USA. 3026899: USENIX Association;  
852 2016. p. 265-83.
- 853 67. Chollet F. *Deep Learning with Python*. Shelter Island, NY: Manning Publications; 2018.
- 854 68. Kingma DP, Ba JL. ADAM: A METHOD FOR STOCHASTIC OPTIMIZATION. 2015.
- 855 69. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, et al. The Ensembl  
856 Variant Effect Predictor. *Genome Biol*. 2016;17(1):122.
- 857 70. Dana JM, Gutmanas A, Tyagi N, Qi G, O'Donovan C, Martin M, et al. SIFTS: updated  
858 Structure Integration with Function, Taxonomy and Sequences resource allows 40-fold increase  
859 in coverage of structure-based annotations for proteins. *Nucleic Acids Res*. 2019;47(D1):D482-  
860 D9.

861