

pubs.acs.org/jcim



Accurately Predicting Mutation-Caused Stability Changes from Protein Sequences Using Extreme Gradient Boosting

3 Xuan Lv, Jianwen Chen, Yutong Lu, Zhiguang Chen, Nong Xiao,* and Yuedong Yang*

Cite This: https://dx.doi.org/10.1021/acs.jcim.0c00064
CCESS Metrics & More Article Recommendations Supporting Information

4 ABSTRACT: Accurately predicting the impact of point mutation on 5 protein stability has crucial roles in protein design and engineering. In this 6 study, we proposed a novel method (BoostDDG) to predict stability 7 changes upon point mutations from protein sequences based on the extreme 8 gradient boosting. We extracted features comprehensively from evolutional 9 information and predicted structures and performed feature selection by a 10 strategy of sequential forward selection. The features and parameters were 11 optimized by homologue-based cross-validation to avoid overfitting. Finally, 12 we found that 14 features from six groups led to the highest Pearson 13 correlation coefficient (PCC) of 0.535, which is consistent with the 0.540 on 14 an independent test. Our method was indicated to consistently outperform 15 other sequence-based methods on three precompiled test sets, and 7363 16 variants on two proteins (PTEN and TPMT). These results highlighted that 17 BoostDDG is a powerful tool for predicting stability changes upon point 18 mutations from protein sequences.

1. INTRODUCTION

19 Nonsynonymous single nucleotide polymorphisms (SNPs) 20 play fundamental roles in protein evolution mainly through 21 inducing changes in amino acid sequences.¹ Previous studies^{2,3} 22 suggest that nonsynonymous SNP mutations can be disease-23 associated and detrimental to human health by affecting 24 structural stabilities of proteins. Even a single residue 25 substitution, namely, point mutation, may have serious effects 26 on a protein's stability and consequently result in changes in 27 the protein's structure and function. Therefore, accurate 28 prediction of mutation-caused protein stability changes 29 ($\Delta\Delta G$) is essential for guiding directed protein evolution 30 and for understanding the relationship between protein 31 variants and diseases.

Over the years, efforts were devoted to developing methods 32 33 to evaluate the stability changes upon point mutations. 34 Traditional mutagenesis studies can accurately quantify the 35 thermodynamic effects of mutations via biological experiments. 36 However, these experimental efforts usually take considerable 37 time and cost, especially for large-scale mutations. With the 38 rapid accumulation of experimentally measured data for both 39 proteins and mutants in the ProTherm database,⁴ computa-40 tional methods were greatly improved in recent years. 41 Predictive stability effects of mutations can be obtained from 42 a protein's sequence or tertiary structure information, and the 43 majority of current methods are structure-based, including 44 several energy-based approaches and machine learning 45 methods. Physical energy functions compute $\Delta\Delta G$ by 46 simulating atomic force fields in a protein structure³ and



may not be applied to large-scale data sets due to high 47 computational requirements. Methods using statistical energy 48 functions calculate the conditional probabilities of certain 49 residues or atoms in different structural environments. For 50 instance, SDM⁶ characterizes free energy changes with the 51 frequency of environment-specific amino acid substitutions, 52 and CUPSAT⁷ utilizes the knowledge-based potentials in a 53 given structure. Empirical potential approaches derive a linear 54 combination of different energy functions, whose weights are 55 fitted to experimental data.^{8,9} Machine learning methods can 56 learn complex nonlinear models from sequence or structural 57 features, and various algorithms have been used to generate 58 predictions, such as support vector machine,¹⁰⁻¹³ decision 59 tree,¹⁴ artificial Neural network,^{15,16} and collaborative filtering 60 model.⁸ 61

One limitation of structure-based approaches is that the 62 tertiary structure of the target protein is often unknown; in 63 fact, less than 0.2% of the proteins in UniProt can retrieve 64 three-dimensional structures in the PDB library.¹⁷ With the 65 development of high-throughput sequencing technologies, the 66 gap between the number of protein sequences and resolved 67

Received: January 21, 2020 Published: March 23, 2020



68 structures is being further widened. In the past, it was generally 69 accepted that the prediction accuracy obtained using sequence 70 alone is lower than using tertiary structure information.¹⁰ 71 Notably, Khan and Vihinen¹⁸ made a systematic analysis of 11 72 kinds of software, indicating that several sequence-based 73 methods^{10,19,20} achieved predictive accuracy comparable to 74 the remaining structural methods. Recent studies have 75 reported that sequence-based predictors^{11,12} achieved similar 76 or better performance than structure-based ones. For the above 77 reasons, we devoted our interest to the sequence-based 78 methods in this work.

Among the machine learning methods used in practice, gradient tree boosting has been shown to give state-of-the-art results in many classification or regression tasks, while rarely being used to predict stability changes upon point mutations. A relatively new ensemble method called extreme gradient boosting (XGBoost) was proposed and shined in many competitions such as KDDCup and Kaggle,²¹ outperforming many other known machine learning methods or deep learning techniques both in accuracy and computational efficiency. Recently published studies have also confirmed the successful applications of XGBoost in computational biology,^{22–24} o including our study to predict genome-wide RNA secondary istructure profiles.²⁵

In this study, we developed a new sequence-based predictor 92 93 (BoostDDG) for protein stability change upon point mutation 94 using the XGBoost machine learning model. We collected and 95 curated a reliable training data set from recently published 96 studies and performed feature selection with parallel model 97 optimization by homologue-based cross-validation. Several 98 studies have pointed out that the predictions of most machine 99 learning methods show a bias toward destabilizing mutations 100 since their training sets are dominated by negative $\Delta\Delta G$ 101 values.^{26–29} To solve this problem, we trained our model using 102 a balanced learning set by introducing reverse mutations with 103 the opposite sign of corresponding $\Delta\Delta G$ values. The selected 104 features in the final input vector were also evaluated through 105 comprehensive ablation experiments. By testing on different 106 data sets, our method was shown to consistently achieve better 107 accuracies in comparison with other state-of-the-art methods. 108 The Web server and all the collected data of BoostDDG are 109 freely available at http://biomed.nscc-gz.cn/server/ 110 BoostDDG/.

2. MATERIALS AND METHODS

2.1. Data Sets. We have collected the training and 111 112 independent test sets from the recent two studies by Folkman 113 et al.¹² and Cao et al.¹⁶ Folkman's study used 1676 mutations 114 for training, and two independent test sets (S236 and S543 115 containing 236 and 543 mutations, respectively). Cao's study 116 used S5444 for training and S276 for testing. We merged 117 S1676 and S5444 data sets and removed mutations whose 118 occurring proteins have a sequence identity of $\geq 25\%$ with any 119 protein in three independent test sets and finally extracted 120 2815 mutations in 150 proteins for training and cross-121 validation. Furthermore, we employed the stability changes 122 on coded protein by PTEN and TPMT genes from the Critical 123 Assessment of Genome Interpretation (CAGI) challenge. After 124 removing the variants with unknown amino acid "X", we 125 independently tested a total of 7363 missense mutations for 126 the PTEN (3736 mutations) and TPMT (3627 mutations) 127 proteins. In this challenge, a stability score was deployed to 128 describe the steady-state abundance of missense mutations

with 0 signifying unstable, 1 signifying wild-type stability, and 129 >1 signifying better stability than the wild-type. This data set 130 was named the CAGI data set. Table S1 summarizes the details 131 of these data sets.

2.2. Candidate Features. In order to train our model, we 133 extracted several classes of candidate features including 134 evolutionary information, predicted structural features, and 135 physicochemical properties for mutant and wild-type amino 136 acids at target residues in the protein sequences. 137

Evolutionary Conservation Features. In this study, we 138 employed the position-specific scoring matrix (PSSM) 139 generated from PSI-BLAST2.7.1³⁰ with an E-value threshold 140 of 0.001 in three iterations to estimate evolutionary 141 conservation of the mutation site. We built three features 142 from the PSSM profile: $\Delta F_{tr} \Delta P_{tr}$ and PE_{tr} . The ΔF_{t} is computed 143 as 144

$$\Delta F_t = (F_{t,mt} - F_{t,wt})/100$$

where $F_{t,mt}$ and $F_{t,wt}$ are the occurrence frequency of mutant 145 (mt) and wild-type (wt) amino acids at the mutation site *t* 146 extracted from the PSSM (the second 20 columns), 147 respectively. In the sites with zeros for all amino acids, the 148 frequencies are extracted from the blosum62 matrix to the 149 wild-type. The ΔP_t is computed as 150

$$\Delta P_t = P_{t,mt} - P_{t,wt}$$

where $P_{t,mt}$ and $P_{t,wt}$ are extracted from the log-likelihood ratio 151 matrix (the first 20 columns) representing the probability of 152 mutant (mt) and wild-type (wt) amino acids at the mutation 153 site *t*. The conservation entropy at mutation site *t* is computed 154 as 155

$$PE_t = \sum_{i=1}^{20} P_{t,i} \cdot F_{t,i}$$

Predicted Structural Features. Since our method is 156 sequence based, we derived structural features according to 157 structural properties predicted from protein sequences, 158 including the relative solvent accessible surface area (rASA), 159 three probabilities (helix, sheet, coil) of the secondary 160 structure (SS), and the disorder probability. The first two 161 classes of features were obtained from SPIDER 3.0,³¹ and the 162 disorder probability was predicted by the SPOT-disorder,³² 163 both with the default parameters. 164

Physicochemical Properties. We employed a group of 165 seven physicochemical parameters provided by Meiler et al.³³ 166 including hydrophobicity, isoelectric point, polarizability, steric 167 parameter, volume, helix tendency, and sheet tendency, 168 namely, AAPh7. we adopted ΔAAPh7 = AAPh7_{mt} – AAPh7_{wt} 169 to calculate the difference between the mutant and wild-type 170 amino acids at the mutation site. We also introduced another 171 feature called AAscore inspired by Liu et al.³⁴ to measure the 172 distinction between wild-type (wt) and mutant (mt) amino 173 acids as 174

$$\theta(R_{mt}, R_{wt}) = \frac{1}{7} \sum_{j}^{\prime} (I_j(R_{mt}) - I_j(R_{wt}))^2$$

where $\theta(R_{mtr}, R_{wt})$ represents the scores calculated using seven 175 physicochemical properties of amino acids presented in the 176 AAindex database (see Table S2). $I_j(R_{mt})$ and $I_j(R_{wt})$ are the 177 normalized values of mutant and wild-type amino acids in 178 property *j*, which can be calculated by 179

Journal of Chemical Information and Modeling

$$I_{j}(a_{i}) = \frac{\tilde{I}_{j}(a_{i}) - \frac{1}{20}\sum_{k=1}^{20}\tilde{I}_{k}(a_{k})}{\sqrt{\frac{1}{20}\sum_{k=1}^{20}(\tilde{I}_{k}(a_{k}) - \frac{1}{20}\sum_{k=1}^{20}\tilde{I}_{k}(a_{k}))^{2}}}$$

¹⁸⁰ where $\tilde{I}_j(a_i)$ symbolizes the original physicochemical parame-¹⁸¹ ters of amino acid a_i in the property j, and $a_k(k = 1,2,3,4\cdots,20)$ ¹⁸² stands for the 20 common amino acids.

2.3. XGBoost Technique. BoostDDG was trained using 184 XGBoost, a novel tree-boosting algorithm to develop an 185 ensemble strong learner from multiple weak learners in an 186 additive way; the tree model outputs a weighted sum of the 187 prediction of each learner by readjusting weights of mislabeled 188 samples in each boosting step. The training procedure can be 189 summarized as follows:

- (i) For each descriptor, sort the feature values and scan thebest splitting point.
- (ii) Select the best splitting point that optimizes theobjective function.
- 194 (iii) Repeat the splitting procedure (the above two steps)
 195 until the predetermined maximum tree depth is
 196 achieved.
- (iv) Assign prediction scores to leaves and prune negativenodes in a bottom-up order.
- (v) Repeat the above steps until the predetermined numberof trees is reached.

201 XGBoost has advantages in preventing overfitting by 202 incorporating regularization terms as well as shrinkage and 203 descriptor subsampling techniques. In addition, distributed and 204 parallel computing enables faster model training.

2.5 2.4. Training Procedure. To find out the optimal set of 206 features and parameters for our model, we coupled feature 207 selection with a parallel grid search procedure. Evaluation of 208 features and parameters in each step was performed using the 209 10-fold cross-validation.

Feature Selection. To ensure the effectiveness and non-210 211 redundancy of all the input features in our model, we selected 212 the best subset of all candidate features using a bidirectional 213 greedy selection algorithm. Both Sequential Forward Selection 214 (SFS) and Sequential Backward Selection (SBS) algorithms 215 were adopted to avoid local optimal results to the greatest 216 extent. SFS (Algorithm 1) begins with an empty set of features 217 F_{sel0} and iteratively searches the best feature f from the 218 remaining feature set F_{left} and adds f to F_{sel0} for a lower root-219 mean-square error (RMŠE). The procedure repeats until F_{left} is 220 empty or RMSE no longer decreases. On the contrary, SBS 221 (Algorithm S1) starts with the full set of features F_{sel0} and 222 removes the worst feature f to yield a lower RMSE. The reverse 223 search process stops when no feature could be discarded. We 224 compared the results generated by the two searching 225 algorithms and selected the one with a higher Pearson's 226 correlation coefficient (PCC) as the final input feature set for 227 our model.

Parallel Grid Search. As the performances are affected by phyperparameters, we have scanned all combinations of parameters by a grid search procedure during the SFS/SBS searches. However, it is considerably time consuming to perform a gird search with the 10-fold cross-validation for each kind of feature set. To speed up the computation, we with implemented an MPI-based approach to parallelize the search process through distributing all the parameter combinations among 288 CPU cores on the Tianhe-2 supercomputer, leading to over 200 times faster than the computation on a single CPU core. Table S3 details the tested hyperparameters 238 set up in the parallel grid search process and the final optimal 239 values. 240

Algorithm 1 Sequential Forward Feature Selection.
Input: The set of all candidate features, F_{list} ;
Output: Features that have been selected, F_{sel0} ;
1: function SFS_Select(F_{list})
2: $F_{left} \leftarrow F_{list}, F_{sel0} \leftarrow \emptyset, loss0 \leftarrow 10$
3: while F_{left} is not empty do
4: $loss_reduction \leftarrow []$
5: for $f \in F_{left}$ do
$6: F_{sel} \leftarrow F_{sel0} \cup f$
7: $bestParams, loss \leftarrow GridSearch(F_{sel})$
8: $loss_reduction.insert(loss0 - loss)$
9: end for
10: if $max(loss_reduction) < 0$ then
11: break
12: else
13: $F_{sel0} \leftarrow F_{sel0} \cup F_{left} [argmaxloss_reduction]$
14: $F_{left} \leftarrow F_{left} - F_{left} [argmaxloss_reduction]$
15: $loss0 \leftarrow loss0 - max(loss_reduction)$
16: end if
17: end while
18: return F_{sel0}
19: end function

Balanced Data Set with Reverse Mutations. Several 241 previous studies have discussed the antisymmetric property 242 of the free energy changes between wild-type and mutant 243 proteins.^{26,28,29,35} The $\Delta\Delta G$ of direct mutation should be 244 equal to the $-\Delta\Delta G$ of reverse variation according to 245 thermodynamic reversibility of mutations. However, most of 246 the predictors ignored the self-consistency requirement due to 247 the bias toward destabilizing mutations in their training sets. 248 Here, we trained our model with a balanced data set by 249 introducing reverse mutations. Each reverse data point has the 250 negative value of experimentally measured $\Delta\Delta G$, and the 251 corresponding features were calculated from the mutated 252 protein sequence. 253

Homologue-Based 10-Fold Cross-Validation. In general 254 cross-validations, the training set can be separated based on 255 mutations (completely random), residues (mutations on the 256 same residue put together), proteins (mutations on the same 257 protein put together), or homologues (mutations on 258 homologous proteins put together). Since mutation, residue, 259 and protein-based cross-validations may share protein 260 information with the training samples, these schemes are likely 261 to cause an overestimate of the performance, as also indicated 262 in the previous study.¹² Relatively, homologue-based evaluation 263 is the strictest assessment since all mutations occurring on the 264 same protein or homologous proteins are always grouped in 265 the same fold to avoid overfitting for a specific group of 266 proteins. We employed the homologue-based scheme on the 267 S2815 data set to conduct our 10-fold cross-validation 268 combined with feature selection and parameter optimization. 269 We strictly separated the proteins from 111 homologous 270 clusters in S2815 into 10 folds to ensure that no two folds 271 share proteins with sequence similarity $\geq 25\%$. The S2815 data 272 set was augmented by introducing hypothetical reverse 273 variations, and each pair of mutations (direct and reverse) 274 was included in the same fold to avoid overfitting. Moreover, 275 we repeated each cross-validation procedure 10 times with 10 276 regenerated 10-folds and averaged the results to gain an 277 unbiased evaluation. 278

Considering that our method is evolutionary based and it is 279 time consuming to generate PSSM profiles from multiple 280 sequence alignment, we also developed another method 281 DeepDDG-single, which is single-sequence based and more 282 283 computationally efficient by removing evolutionary conserva-284 tion information and replacing the predicted secondary 285 structure features from SPIDER 3.0³¹ with those generated by 286 SPIDER3-single³⁶ that also made predictions from only 287 sequences without using evolutionary information.

3. RESULTS

3.1. Feature Selection. We applied a bidirectional greedy feature selection algorithm on the S2815 data set with 10-fold

Table 1. PCC between Experimental $\Delta\Delta G$ and $\Delta\Delta G$ Predicted by Using SFS-Selected Feature Groups or by Removing Each Feature Group from the Final Model

Feature groups ^a	CV ^c	Ind. test ^c	Feature groups ^b	CV ^c	Ind. test ^c
$\Delta AAPh7$	0.385	0.332	BoostDDG	0.535	0.540
+rASA	0.461	0.440	$-\Delta AAPh7$	0.325	0.451
$+\Delta F_t$	0.510	0.508	-rASA	0.502	0.511
+SS	0.524	0.526	-AAscore	0.529	0.537
+AAscore	0.531	0.538	$-\Delta F_t$	0.478	0.471
+disorder	0.535	0.540	-disorder	0.523	0.544
all features	0.530	0.531	-SS	0.520	0.517
SVM ^d	0.515	0.505	SVM ^e	0.519	0.511

^{*a*}Additive feature groups selected in SFS algorithm. ^{*b*}Removed feature group from the final model. ^{*c*}CV, 10-fold cross-validation; Ind. test, independent test on S543. ^{*d*}Performance of SVM model using all features. ^{*c*}Performance of SVM model using selected features.

Table 2. Comparison of BoostDDG and Related Methods on Three Benchmark Data $Sets^{a}$

	S236 test set		S543 test set		S276 test set	
Methods	PCC	RMSE	PCC	RMSE	PCC	MAE
I Mutant ^b	0.443	1.18	0.323	1.37	0.391	1.08
I Mutant ^c	0.521	1.07	0.356	1.34	0.453	0.91
MUpro	0.362	1.20	0.332	1.32	0.190	1.06
Rosetta	0.270	1.88	0.380	3.58	0.339	5.25
FoldX	0.277	1.70	0.405	1.87	0.300	2.13
DFIRE	0.535	1.18	0.450	1.44	0.230	1.25
PoPMuSiC	0.570	1.05	0.533	1.21	0.443	0.91
EASE-MM	0.589	1.03	0.530	1.22	0.402	0.91
STRUM	-	-	-	_	0.447	0.88
mCSM	_	-	_	_	0.467	0.90
INPS ^d	-	-	_	-	0.474	0.89
SDM	-	-	-	_	0.483	1.02
BoostDDG-single	0.511	1.17	0.431	1.39	0.355	1.02
BoostDDG	0.600	1.01	0.540	1.21	0.514	0.78

^{*a*}Predictions of S236 and S543 were calculated from the data provided in Folkman's study.¹² Predictions of S276 were collected from Cao's study.¹⁶. ^{*b*}Sequence-based I Mutant. ^{*c*}Structure-based I Mutant. ^{*d*}Computed by ourselves.

290 cross-validation. The final combination of predictive features 291 for BoostDDG was composed of the following features: 292 Δ AAPh7 (the changes of hydrophobicity, volume, polar-293 izability, isoelectric, helix tendency, sheet tendency, and steric), 294 rASA, ΔF_p , AAscore, SS (helix probability, sheet probability, 295 and coil probability), and disorder probability. All these 296 features were real-valued and constituted a 14-dimensional 297 input vector. Table 1 (left) displays the performances of 298 XGBoost-based models using additive feature groups in each 299 iteration of the SFS algorithm. The SFS selected six feature

t1

pubs.acs.org/jcim

Table 3. Independent Test Results for PTEN and TPMT Data Sets in Terms of PCC and RMSE between Expected Stability Scores and Linear-Fitted Predictions

	PTEN		TPMT		TPMT and PTEN	
Methods	PCC	RMSE	PCC	RMSE	PCC	RMSE
MUpro1.1	0.221	0.320	0.235	0.349	0.229	0.336
I Mutant2.0	0.199	0.322	0.251	0.348	0.228	0.336
EASE-MM	0.410	0.299	0.379	0.332	0.391	0.317
STRUM	0.134	0.325	0.426	0.325	0.320	0.327
INPS	0.444	0.294	0.383	0.332	0.418	0.313
BoostDDG	0.464	0.290	0.420	0.323	0.440	0.310

groups containing totally 14 features. The trained model 300 achieved a PCC of 0.535, slightly higher than the 0.530 by 301 using all features. The robustness of our trained model with 302 selected features further proved essentially same results on the 303 independent test of S543. The PCC of 0.540 is consistently 304 higher than 0.531 by using all features. The results suggest that 305 our feature selection procedure is helpful to improve 306 prediction performance by removing redundancy and 307 irrelevant features. 308

For comparison, we also trained SVM regression models 309 implemented in the scikit-learn³⁷ python library using all 310 features and selected features based on the S2815 data set. 311 Both SVM models achieved PCCs lower than those generated 312 by BoostDDG whether in cross-validation (0.515 and 0.519) 313 or the independent test (0.505 and 0.511), indicating that 314 XGBoost is better for this task. We also tested the SBS 315 algorithm to remove redundant features from the whole feature 316 set, but no feature was filtered out. 317

3.2. Assessment of Feature Importance. To get further 318 insights into the factors that contribute to the mutation- 319 induced protein stability changes, we employed ablation 320 experiments by removing each feature group from the final 321 model and re-evaluating the performance based on the cross- 322 validation and test set (Table 1, right). Both the PCCs of 323 cross-validation and the independent test decrease the most 324 when removing Δ AAPh7. This is understandable since 325 Δ AAPh7 is a seven-dimensional feature vector that is larger 326 than other feature groups. Among the remaining features, the 327 exclusion of ΔF_t caused the largest drop of PCC, confirming 328 that evolutionary conservation is an essential factor in the 329 prediction of stability changes. The predicted solvent 330 accessibility has a significant effect on our model as well, 331 which is consistent with the results by Dehghanpoor et al.³⁸ 332 and Cao et al.,¹⁶ likely because residues with smaller accessible 333 surface areas (buried) are more vulnerable to destabilizing 334 substitutions. Different from the positive contribution of 335 disorder probability in the cross-validation, the removal of 336 disorder probability led to a slight increase in PCC from 0.540 337 to 0.544 in the test set, likely due to the difference of data sets, 338 so we still kept the disorder probability in our final model. 339

3.3. Method Comparison on Three Benchmark Data ³⁴⁰ **Sets.** We adopted three different benchmark test sets, S236 ³⁴¹ together with S543 from Folkman et al.¹² and S276 from Cao ³⁴² et al.¹⁶ to verify if BoostDDG generalizes well. None of the ³⁴³ three test sets was used for training or optimizing our method. ³⁴⁴ For comparing with other predictors, we only tested the ³⁴⁵ experimentally measured $\Delta\Delta G$ values (direct mutations). We ³⁴⁶ compared our method with 11 state-of-the-art methods, ³⁴⁷ including four sequence-based methods (I Mutant, ¹⁹ ³⁴⁸ MUpro, ¹⁰ INPS, ¹¹ and EASE-MM¹²) and seven structure- ³⁴⁹



Figure 1. Comparison between expected stability scores and scores predicted with six sequence-based methods for the CAGI data set.

350 based energy functions (Rosetta,³⁹ FoldX,⁴⁰ DFIRE,⁴¹ 351 PoPMuSiC,⁴² STRUM,¹⁷ mCSM,⁴³ SDM⁶). The structure-352 based version of I Mutant was also included. We tested a few 353 methods (STRUM, mCSM, INPS, and SDM) only on the 354 S276 set since their training sets contain many proteins 355 sequentially similar to S543 and S236 sets.

The performance was assessed in terms of PCC and RMSE. 356 357 As S276 does not have reported RMSE in the previous study,¹⁶ we used the reported mean absolute error (MAE) instead on 358 this data set for a consistent comparison. Table 2 outlines the 359 results obtained with each of the above-mentioned methods on 360 three test sets. Remarkably, BoostDDG produced the highest 361 PCC and lowest RMSE (or MAE) for all three data sets, 362 outperforming all sequence-based and structure-based meth-363 ods tested in previous studies. BoostDDG achieved a higher 364 365 PCC in the S236 but a lower PCC in the S276 than the one in 366 S543. We noticed that the PCC in the S236 is even higher than 367 the PCC (= 0.54) in our training set S2815. This is likely 368 because S2815 and S276 data sets contain $\Delta\Delta G$ values 369 obtained under various experimental conditions, while the 370 S236 and S543 data sets only contain the $\Delta\Delta G$ measurements under the standard pH and temperature. Thus, S236 is less 371 challenging as also seen by a typically higher PCC in S236 and 372 a lower PCC in S276 by other methods as well. Figures S1-S3373 show the predicted versus experimentally measured $\Delta\Delta G$ on 374 375 three test sets, respectively.

t2

We also tested our single-sequence-based method 376 BoostDDG-single on these data sets. Despite the gap with 377 evolutionary-based methods, it still surpasses some popular 378 sequence-based predictors (I Mutant¹⁹ and MUpro¹⁰). In fact, 379 the availability of evolutionary information can be limited since 380 most proteins have few known homologous sequences.³⁶ 381 382 Additionally, the generation of evolutionary sequence profile is 383 time consuming even for a short sequence. Therefore, the 384 single-sequence-based method is able to constitute a trade-off 385 between accuracy and computational efficiency.

3.4. Independent Testing of CAGI Data Set. We further 386 tested our method on two proteins, PTEN (phosphatase and 387 TEnsin homologue) and TPMT (thiopurine S-methyl trans- 388 ferase), from the CAGI challenge that was provided by Fowler 389 and Fields.⁴⁴ Both proteins have pairwise identities of less than 390 25% with proteins in our training set. We used the full-length 391 sequence of the two proteins to compare BoostDDG with five 392 other leading sequence-based approaches (MUpro1.1,¹⁰ I 393 EASE-MM,¹² STRUM,¹⁷ and INPS¹¹). 394 Mutant2.0,¹⁹ STRUM is a structure-based method but can be applied to 395 protein sequences with predicted 3D structures. Pejaver et al. 396 analyzed the performances of several computational predictors 397 on the PTEN and TPMT protein variants and found that 398 predictors ranked different by metrics.⁴⁵ For consistent 399 comparisons, we calculated the PCC and RMSE between 400 stability scores and the predictions of the six methods.

Table 3 outlines the PCC between stability scores and the 402 t3 prediction data of the six methods. BoostDDG produced a 403 more balanced and accurate prediction for the two proteins 404 with a PCC of 0.420-0.464. The second best method is INPS, 405 achieving a PCC of 0.44 in total, which is significantly lower 406 than our method according to the Fisher test (P = 0.05). 407 Although STRUM yielded a comparable correlation (0.426) 408 for the TPMT data, the accuracy for PTEN data (0.134) is 409 much lower than other methods probably because the accuracy 410 of the simulated 3D structure decreases when the sequence 411 length is too long (>400). The independent test further 412 demonstrates that our method is able to give an improved 413 prediction when experimental structures are not available. 414 Figure 1 illustrates the expected stability score versus predicted 415 fl $\Delta \Delta G.$ 416

4. **DISCUSSION**

In this work, we have developed a new sequence-based 417 method, BoostDDG, to predict the stability effects of point 418 mutations based on XGBoost. We carefully selected features 419

420 from initial input with a bidirectional greedy algorithm to 421 generate the best nonredundant feature set. To select features 422 in a fair and efficient way, we performed parallel model 423 optimization in each step of feature selection with the strictest 424 homologue-based 10-fold cross-validation, in which no two 425 folds shared homologous proteins.

We employed three different independent test sets from 426 427 Folkman et al.¹² and Cao et al.¹⁶ and made a comparison with 428 other methods tested in their works. The results showed that 429 BoostDDG produced consistently the best predictions 430 measured in both PCC and RMSE. Furthermore, we 431 independently tested our method with the PTEN and 432 TPMT data sets from the CAGI challenge and found that 433 BoostDDG achieved a balanced correlation of 0.42-0.46 that 434 is superior to those generated by other sequence-based 435 methods. This can be largely attributed to our strict 436 homologue-based validation scheme and the robust regulariza-437 tion performance of XGBoost. Since a majority of the previous 438 methods were trained and evaluated using a residue-based 439 scheme or a protein-based scheme with a high sequence 440 identity cutoff, their performance may be overestimated due to 441 the correlation between available data sets. XGBoost is a 442 regularized scalable extension of traditional tree-boosting 443 algorithms which is much less susceptible to overfitting by 444 tuning the corresponding parameters. Following the previous 445 studies,^{11,16} we included the thermodynamic reversibility of 446 variations to enhance the symmetry of our model. We found 447 the inclusion slightly increased the PCC on the CV (0.530 to 448 0.535) relative to the one without the inclusion and led to 449 small and divergent changes on four test sets (0.64-0.60 on 450 S236, 0.538-0.54 on S543, 0.49-0.51 on S276, and 0.45-0.44 451 on the CAGI data sets). Although far from extensive, the 452 independent test performance highlights that as a sequence-453 based method BoostDDG is a promising alternative to existing 454 popular predictors when experimental structures are not 455 available.

Typically, the stability change $\Delta\Delta G$ is sensitive to 456 457 experimental conditions such as pH and temperature. 458 However, these environmental parameters were not encoded 459 into input features considering that the number of measure-460 ments of each mutation at different environmental conditions 461 is too small to provide generalized predictions. Despite the 462 favorable independent test performance that our method has 463 achieved, the correlation values became lower on the larger test 464 sets, indicating that further improvement should be considered. 465 Given that our method only focuses on the features extracted 466 from the mutation site, we aim to design more useful features 467 to capture the interaction changes between neighboring 468 residues caused by mutations in our future work.

ASSOCIATED CONTENT 469

470 Supporting Information

471 The Supporting Information is available free of charge at 472 https://pubs.acs.org/doi/10.1021/acs.jcim.0c00064.

473 Table S1: Summary of data sets used. Table S2: Seven physicochemical properties of amino acids used for 474 AAscore. Table S3: Hyperparameters set up in the 475 parallel grid search process and the optimal values used 476 for the final model. Algorithm S1: Sequential Backward 477 Selection algorithm. Figure S1: Experimentally measured 478 $\Delta\Delta G$ from the S236 data set versus $\Delta\Delta G$ predicted with 479 the four sequence-based methods (BoostDDG, I 480

pubs.acs.org/jcim

491

Mutant2.0, MUpro1.1, and EASE-MM) and four 481 structure-based methods (I Mutant2.0, FoldX, DFIRE, 482 and PoPMuSiC2.1). Figure S2: Experimentally meas- 483 ured $\Delta\Delta G$ from the S543 data set versus $\Delta\Delta G$ predicted 484 with the four sequence-based methods (BoostDDG, I 485 Mutant2.0, MUpro1.1, and EASE-MM) and four 486 structure-based methods (I Mutant2.0, FoldX, DFIRE, 487 and PoPMuSiC2.1). Figure S3: Experimentally meas- 488 ured $\Delta\Delta G$ from the S276 data set versus $\Delta\Delta G$ obtained 489 from BoostDDG. (PDF) 490

AUTHOR INFORMATION

Corresponding Authors	492
Yuedong Yang – School of Data and Computer Science, Sun	493
Yat-sen University, Guangzhou, Guangdong 510275, China;	494
Key Laboratory of Machine Intelligence and Advanced	495
Computing, Sun Yat-sen University, Ministry of Education,	496
Guangzhou, Guangdong 510275, China; 💿 orcid.org/0000-	497
0002-6782-2813; Email: yangyd25@mail.sysu.edu.cn	498
Nong Xiao – State Key Laboratory of High-Performance	499
Computing, School of Computer Science, National University of	500
Defense Technology, Changsha, Hunan 410073, China; School	501
of Data and Computer Science, Sun Yat-sen University,	502
Guangzhou, Guangdong 510275, China; Email: xiaon6@	503
mail.sysu.edu.cn	504
Authors	505
Authors	505
Authors Xuan Lv – State Key Laboratory of High-Performance	505 506
Authors Xuan Lv – State Key Laboratory of High-Performance Computing, School of Computer Science, National University of Defense Technology Chargebra Human 410072 Ching	505 506 507
Authors Xuan Lv – State Key Laboratory of High-Performance Computing, School of Computer Science, National University of Defense Technology, Changsha, Hunan 410073, China;	505 506 507 508
Authors Xuan Lv − State Key Laboratory of High-Performance Computing, School of Computer Science, National University of Defense Technology, Changsha, Hunan 410073, China; © orcid.org/0000-0001-8821-0026	505 506 507 508 509
 Authors Xuan Lv – State Key Laboratory of High-Performance Computing, School of Computer Science, National University of Defense Technology, Changsha, Hunan 410073, China; orcid.org/0000-0001-8821-0026 Jianwen Chen – School of Data and Computer Science, Sun 	505 506 507 508 509 510
 Authors Xuan Lv – State Key Laboratory of High-Performance Computing, School of Computer Science, National University of Defense Technology, Changsha, Hunan 410073, China; orcid.org/0000-0001-8821-0026 Jianwen Chen – School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong 510275, China; 	505 506 507 508 509 510 511
 Authors Xuan Lv – State Key Laboratory of High-Performance Computing, School of Computer Science, National University of Defense Technology, Changsha, Hunan 410073, China; orcid.org/0000-0001-8821-0026 Jianwen Chen – School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong 510275, China; orcid.org/0000-0001-7999-2070 	505 506 507 508 509 510 511 512
 Authors Xuan Lv – State Key Laboratory of High-Performance Computing, School of Computer Science, National University of Defense Technology, Changsha, Hunan 410073, China; orcid.org/0000-0001-8821-0026 Jianwen Chen – School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong 510275, China; orcid.org/0000-0001-7999-2070 Yutong Lu – School of Data and Computer Science, Sun Yat-sen 	505 506 507 508 509 510 511 512 513
 Authors Xuan Lv – State Key Laboratory of High-Performance Computing, School of Computer Science, National University of Defense Technology, Changsha, Hunan 410073, China; orcid.org/0000-0001-8821-0026 Jianwen Chen – School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong 510275, China; orcid.org/0000-0001-7999-2070 Yutong Lu – School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong 510275, China 	505 506 507 508 509 510 511 512 513 514
 Authors Xuan Lv – State Key Laboratory of High-Performance Computing, School of Computer Science, National University of Defense Technology, Changsha, Hunan 410073, China; orcid.org/0000-0001-8821-0026 Jianwen Chen – School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong 510275, China; orcid.org/0000-0001-7999-2070 Yutong Lu – School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong 510275, China Zhiguang Chen – School of Data and Computer Science, Sun 	505 506 507 508 509 510 511 512 513 514 515
 Authors Xuan Lv – State Key Laboratory of High-Performance Computing, School of Computer Science, National University of Defense Technology, Changsha, Hunan 410073, China; orcid.org/0000-0001-8821-0026 Jianwen Chen – School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong 510275, China; orcid.org/0000-0001-7999-2070 Yutong Lu – School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong 510275, China Zhiguang Chen – School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong 510275, China 	505 506 507 508 509 510 511 512 513 514 515 516

Con https://pubs.acs.org/10.1021/acs.jcim.0c00064 518

Notes

The authors declare no competing financial interest. 520

ACKNOWLEDGMENTS

This work was supported by the National Key R&D Program 522 of China (2018YFC0910500), National Natural Science 523 Foundation of China (U1611261, 61772566, and 81801132), 524 Guangdong Frontier & Key Tech Innovation Program 525 (2018B010109006, 2019B020228001), and Introducing In- 526 novative and Entrepreneurial Teams (2016ZT06D211). We 527 appreciate the National Supercomputer Center in Guangzhou 528 for providing us with computing resources. 529

ABBREVIATIONS

530

519

521

PCC, Pearson correlation coefficient; XGBoost, extreme 531 gradient boosting; SS, secondary structure; rASA, relative 532 solvent accessible surface area; SNPs, single nucleotide 533 polymorphisms; PSSM, position-specific scoring matrix; SFS, 534 Sequential Forward Selection; SBS, Sequential Backward 535 Selection; CAGI, critical assessment of genome interpretation 536

pubs.acs.org/jcim

537 **REFERENCES**

(1) Tennessen, J. A.; Bigham, A. W.; O'Connor, T. D.; Fu, W.;
Kenny, E. E.; Gravel, S.; McGee, S.; Do, R.; Liu, X.; Jun, G.; et al.
Evolution and functional impact of rare coding variation from deep
sequencing of human exomes. *Science* 2012, 337, 64–69.

542 (2) Tokuriki, N.; Tawfik, D. S. Stability effects of mutations and 543 protein evolvability. *Curr. Opin. Struct. Biol.* **2009**, *19*, 596–604.

544 (3) Yates, C. M.; Sternberg, M. J. The effects of non-synonymous 545 single nucleotide polymorphisms (nsSNPs) on protein-protein 546 interactions. *J. Mol. Biol.* **2013**, 425, 3949–63.

547 (4) Kumar, M. D.; Bava, K. A.; Gromiha, M. M.; Prabakaran, P.; 548 Kitajima, K.; Uedaira, H.; Sarai, A. ProTherm and ProNIT: 549 thermodynamic databases for proteins and protein-nucleic acid 550 interactions. *Nucleic Acids Res.* **2006**, *34*, D204–6.

551 (5) Benedix, A.; Becker, C. M.; de Groot, B. L.; Caflisch, A.; 552 Böckmann, R. A. Predicting free energy changes using structural 553 ensembles. *Nat. Methods* **2009**, *6*, 3.

(6) Pandurangan, A. P.; Ochoa-Montano, B.; Ascher, D. B.; Blundell,
555 T. L. SDM: a server for predicting effects of mutations on protein
556 stability. *Nucleic Acids Res.* 2017, 45, W229–W235.

(7) Parthiban, V.; Gromiha, M. M.; Schomburg, D. CUPSAT:
prediction of protein stability upon point mutations. *Nucleic Acids Res.*2006, 34, W239–W242.

(8) Wainreb, G.; Wolf, L.; Ashkenazy, H.; Dehouck, Y.; Ben-Tal, N.
561 Protein stability: a single recorded mutation aids in predicting the
562 effects of other mutations in the same amino acid site. *Bioinformatics*563 2011, 27, 3286-3292.

564 (9) Schymkowitz, J.; Borg, J.; Stricher, F.; Nys, R.; Rousseau, F.; 565 Serrano, L. The FoldX web server: an online force field. *Nucleic Acids* 566 *Res.* **2005**, 33, W382–W388.

567 (10) Cheng, J.; Randall, A.; Baldi, P. Prediction of protein stability 568 changes for single-site mutations using support vector machines. 569 *Proteins: Struct., Funct., Genet.* **2006**, *62*, 1125–1132.

570 (11) Fariselli, P.; Martelli, P. L.; Savojardo, C.; Casadio, R. INPS: 571 predicting the impact of non-synonymous variations on protein 572 stability from sequence. *Bioinformatics* **2015**, *31*, 2816–2821.

573 (12) Folkman, L.; Stantic, B.; Sattar, A.; Zhou, Y. EASE-MM: 574 Sequence-based prediction of mutation-induced stability changes with 575 feature-based multiple models. *J. Mol. Biol.* **2016**, *428*, 1394–1405.

576 (13) Capriotti, E.; Fariselli, P.; Calabrese, R.; Casadio, R. Predicting 577 protein stability changes from sequences using support vector 578 machines. *Bioinformatics* **2005**, *21*, ii54–ii58.

579 (14) Huang, L.-T.; Gromiha, M. M.; Ho, S.-Y. iPTREE-STAB: 580 interpretable decision tree based method for predicting protein 581 stability changes upon mutations. *Bioinformatics* **2007**, *23*, 1292– 582 1293.

583 (15) Capriotti, E.; Fariselli, P.; Casadio, R. A neural-network-based 584 method for predicting protein stability changes upon single point 585 mutations. *Bioinformatics* **2004**, *20*, i63–i68.

(16) Cao, H.; Wang, J.; He, L.; Qi, Y.; Zhang, J. Z. DeepDDG:
Predicting the Stability Change of Protein Point Mutations Using
Neural Networks. J. Chem. Inf. Model. 2019, 59, 1508–1514.

S89 (17) Quan, L.; Lv, Q.; Zhang, Y. STRUM: structure-based
S90 prediction of protein stability changes upon single-point mutation.
S91 *Bioinformatics* 2016, 32, 2936–2946.

592 (18) Khan, S.; Vihinen, M. Performance of protein stability 593 predictors. *Hum. Mutat.* **2010**, *31*, 675–684.

594 (19) Capriotti, E.; Fariselli, P.; Casadio, R. I-Mutant2. 0: predicting 595 stability changes upon mutation from the protein sequence or 596 structure. *Nucleic Acids Res.* **2005**, 33, W306–W310.

597 (20) Dosztányi, Z.; Fiser, A.; Simon, I. Stabilization centers in
598 proteins: identification, characterization and predictions. *J. Mol. Biol.*599 **1997**, 272, 597–612.

600 (21) Chen, T.; Guestrin, C. Xgboost: A scalable tree boosting 601 system. In *Proceedings of the 22nd ACM SIGKDD International* 602 *Conference on Knowledge Discovery and Data Mining*, San Francisco, 603 August 2016; pp 785–794.

604 (22) Mendik, P.; Dobronyi, L.; Hári, F.; Kerepesi, C.; Maia-Moço, 605 L.; Buszlai, D.; Csermely, P.; Veres, D. V. Translocatome: a novel resource for the analysis of protein translocation between cellular 606 organelles. *Nucleic Acids Res.* **2019**, 47, D495–D505. 607

(23) Zou, L. S.; Erdos, M. R.; Taylor, D. L.; Chines, P. S.; Varshney, 608 A.; Parker, S. C.; Collins, F. S.; Didion, J. P. BoostMe accurately 609 predicts DNA methylation values in whole-genome bisulfite 610 sequencing of multiple human tissues. *BMC Genomics* **2018**, *19*, 390. 611 (24) Chen, X.; Huang, L.; Xie, D.; Zhao, Q. EGBMMDA: extreme 612

gradient boosting machine for MiRNA-disease association prediction. 613 Cell Death Dis. 2018, 9, 3.

(25) Ke, Y.; Rao, J.; Zhao, H.; Lu, Y.; Xiao, N.; Yang, Y. Accurate 615 Prediction of Genome-wide RNA Secondary Structure Profile Based 616 On Extreme Gradient Boosting. *bioRxiv*, April 16, **2019**. 617 DOI: 10.1101/610782. 618

(26) Pucci, F.; Bernaerts, K. V.; Kwasigroch, J. M.; Rooman, M. 619 Quantification of biases in predictions of protein stability changes 620 upon mutations. *Bioinformatics* **2018**, *34*, 3659–3665. 621

(27) Usmanova, D. R; Bogatyreva, N. S; Arino Bernad, J.; Eremina, 622 A. A; Gorshkova, A. A; Kanevskiy, G. M; Lonishin, L. R; Meister, A. 623 V; Yakupova, A. G; Kondrashov, F. A; Ivankov, D. N Self-consistency 624 test reveals systematic bias in programs for prediction change of 625 stability upon mutation. *Bioinformatics* **2018**, *34*, 3653–3658. 626

(28) Fang, J. A Critical Review of Five Machine Learning-Based 627 Algorithms for Predicting Protein Stability Changes upon Mutation. 628 *Briefings Bioinf.* **2019**, *na*, na DOI: 10.1093/bib/bbz071. 629

(29) Thiltgen, G.; Goldstein, R. A. Assessing predictors of changes 630 in protein stability upon mutation using self-consistency. *PLoS One* 631 **2012**, 7, e46084. 632

(30) Schäffer, A. A.; Aravind, L.; Madden, T. L.; Shavirin, S.; 633 Spouge, J. L.; Wolf, Y. I.; Koonin, E. V.; Altschul, S. F. Improving the 634 accuracy of PSI-BLAST protein database searches with composition- 635 based statistics and other refinements. *Nucleic acids research* 2001, 29, 636 2994–3005. 637

(31) Heffernan, R.; Dehzangi, A.; Lyons, J.; Paliwal, K.; Sharma, A.; 638 Wang, J.; Sattar, A.; Zhou, Y.; Yang, Y. Highly accurate sequence- 639 based prediction of half-sphere exposures of amino acid residues in 640 proteins. *Bioinformatics* **2016**, *32*, 843–849. 641

(32) Hanson, J.; Yang, Y.; Paliwal, K.; Zhou, Y. Improving protein 642 disorder prediction by deep bidirectional long short-term memory 643 recurrent neural networks. *Bioinformatics* **2016**, *33*, 685–692. 644

(33) Meiler, J.; Müller, M.; Zeidler, A.; Schmäschke, F. Generation 645 and evaluation of dimension-reduced amino acid parameter 646 representations by artificial neural networks. *J. Mol. Model.* **2001**, *7*, 647 360–369. 648

(34) Liu, B.; Wang, S.; Wang, X. DNA binding protein identification 649 by combining pseudo amino acid composition and profile-based 650 protein representation. *Sci. Rep.* **2015**, *5*, 15479. 651

(35) Capriotti, E.; Fariselli, P.; Rossi, I.; Casadio, R. A three-state 652 prediction of single point mutations on protein stability changes. 653 *BMC Bioinf.* **2008**, *9*, S6. 654

(36) Heffernan, R.; Paliwal, K.; Lyons, J.; Singh, J.; Yang, Y.; Zhou, 655 Y. Single-sequence-based prediction of protein secondary structures 656 and solvent accessibility by deep whole-sequence learning. *J. Comput.* 657 *Chem.* **2018**, 39, 2210–2216. 658

(37) Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; 659 Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; 660 Dubourg, V. Scikit-learn: Machine learning in Python. *J. Mach. Learn.* 661 *Res.* **2011**, *12*, 2825–2830. 662

(38) Dehghanpoor, R.; Ricks, E.; Hursh, K.; Gunderson, S.; 663 Farhoodi, R.; Haspel, N.; Hutchinson, B.; Jagodzinski, F. Predicting 664 the effect of single and multiple mutations on protein structural 665 stability. *Molecules* **2018**, 23, 251. 666

(39) Kellogg, E. H.; Leaver-Fay, A.; Baker, D. Role of conforma- 667 tional sampling in computing mutation-induced changes in protein 668 structure and stability. *Proteins: Struct., Funct., Genet.* **2011**, *79*, 830– 669 838. 670

(40) Guerois, R.; Nielsen, J. E.; Serrano, L. Predicting changes in the 671 stability of proteins and protein complexes: a study of more than 1000 672 mutations. *J. Mol. Biol.* **2002**, *320*, 369–387. 673

674 (41) Zhou, H.; Zhou, Y. Distance-scaled, finite ideal-gas reference 675 state improves structure-derived potentials of mean force for structure 676 selection and stability prediction. *Protein science* **2002**, *11*, 2714– 677 2726.

678 (42) Dehouck, Y.; Grosfils, A.; Folch, B.; Gilis, D.; Bogaerts, P.; 679 Rooman, M. Fast and accurate predictions of protein stability changes 680 upon mutations using statistical potentials and neural networks: 681 PoPMuSiC-2.0. *Bioinformatics* **2009**, *25*, 2537–2543.

(43) Pires, D. E.; Ascher, D. B.; Blundell, T. L. mCSM: predicting
683 the effects of mutations in proteins using graph-based signatures.
684 Bioinformatics 2014, 30, 335-342.

685 (44) Fowler, D. M.; Fields, S. Deep mutational scanning: a new style 686 of protein science. *Nat. Methods* **2014**, *11*, 801.

687 (45) Pejaver, V.; Babbi, G.; Casadio, R.; Folkman, L.; Katsonis, P.;

688 Kundu, K.; Lichtarge, O.; Martelli, P. L.; Miller, M.; Moult, J.; et al.

689 Assessment of methods for predicting the effects of PTEN and TPMT 690 protein variants. *Hum. Mutat.* **2019**, *40*, 1495–1506.