

22 **Abstract:**

23 There is strong interest in accurate methods for predicting changes in protein stability resulting
24 from amino acid mutations to the protein sequence. Recombinant proteins must often be stabilized
25 to be used as therapeutics or reagents, and destabilizing mutations are implicated in a variety of
26 diseases. Due to increased data availability and improved modeling techniques, recent studies have
27 shown advancements in predicting changes in protein stability when a single point mutation is
28 made. Less focus has been directed toward predicting changes in protein stability when there are
29 two or more mutations, despite the significance of mutation clusters for disease pathways and
30 protein design studies. Here, we analyze the largest available dataset of double point mutation
31 stability and benchmark several widely used protein stability models on this and other datasets.
32 We identify a blind spot in how predictors are typically evaluated on multiple mutations, finding
33 that, contrary to assumptions in the field, current stability models are unable to consistently capture
34 epistatic interactions between double mutations. We observe one notable deviation from this trend,
35 which is that epistasis-aware models provide marginally better predictions on stabilizing double
36 point mutations. We develop an extension of the ThermoMPNN framework for double mutant
37 modeling as well as a novel data augmentation scheme which mitigates some of the limitations in
38 available datasets. Collectively, our findings indicate that current protein stability models fail to
39 capture the nuanced epistatic interactions between concurrent mutations due to several factors,
40 including training dataset limitations and insufficient model sensitivity.

41

42 **Keywords:**

43 Protein stability, epistasis, point mutations, deep learning, protein design

44

45 **Significance:**

46 Protein stability is governed in part by epistatic interactions between energetically coupled

47 residues. Prediction of these couplings represents the next frontier in protein stability modeling. In

48 this work, we benchmark protein stability models on a large dataset of double point mutations and

49 identify previously overlooked limitations in model design and evaluation. We also introduce

50 several new strategies to improve modeling of epistatic couplings between protein point mutations.

51

52 **Introduction:**

53 Thermodynamic stability is an important property that can impact the fitness of a protein^{1,2}.
54 Molecular biologists often introduce mutations to probe structure-function relationships within
55 proteins, and aberrant stability profiles are implicated in a variety of diseases^{3,4}. Additionally, as
56 engineered proteins are increasingly used as therapeutics⁵ and research tools⁶, their stability must
57 often be optimized to improve production yields and efficacy⁷. Recent advancements in assay
58 design and transfer learning have enabled deep neural networks to predict the change in stability
59 ($\Delta\Delta G$) caused by single point mutations faster and more accurately than prior approaches⁸⁻¹⁰.
60 However, relatively few studies have attempted to explicitly model multiple point mutations, for
61 a few reasons. Not only is reliable stability data less abundant for multiple mutations, but the
62 possible mutation space also increases exponentially with the number of mutations, resulting in a
63 sparse energy landscape that is difficult to model.

64 In this study, we focus on the task of predicting changes in stability ($\Delta\Delta G$) caused by double
65 point mutations. We are partially motivated by the observation that protein engineers often seek to
66 identify clusters of two or more mutations which may improve stability beyond levels achievable
67 with single mutant sweeps through favorable couplings such as hydrogen bonding or apolar
68 packing¹¹. Also, biological researchers must sometimes contend with multiple concurrent
69 mutations introduced by cancer¹², bacteria¹³, or viruses¹⁴. A single mutant stability model can be
70 used to approximate double mutant $\Delta\Delta G$ by simply adding the two constituent single mutant
71 contributions. The drawback of this additive approach is that it omits any epistatic coupling that
72 may arise from the interaction of the two mutations. As such, the utility of a double mutant
73 predictor is derived from its ability to provide improvements relative to the additive predictions
74 provided from its equivalent single mutant model. Despite this observation, double mutant stability

75 models are rarely evaluated in this way. We posit that this represents a significant blind spot in our
76 current understanding of protein stability models, which we aim to address in this work.

77 To that end, we develop a novel method for modeling stability changes due to double point
78 mutations which we call ThermoMPNN-Double (“ThermoMPNN-D”). We analyze the largest
79 available double mutant dataset and introduce a new data augmentation protocol to address
80 shortcomings in data availability. We evaluate ThermoMPNN-D against popular methods from the
81 literature, and we take the additional step of evaluating each predictor against its own additive
82 equivalent. We show that ThermoMPNN-D and its single mutant analogue, ThermoMPNN,
83 provide competitive performance on two datasets of double mutants gathered on a diverse set of
84 proteins. We use deep mutational scanning (DMS) data as an orthogonal test set, finding that the
85 methods Mutate Everything and FoldX perform the best on this task. Overall, we find that
86 epistasis-aware double mutant models rarely outperform their single mutant counterparts, with the
87 notable exception that they provide improved prediction of stabilizing double mutants.

88

89 **Results:**

90 *Adapting ThermoMPNN to model double mutations*

91 We developed a novel neural network, ThermoMPNN-D, to model double point mutations by
92 making several modifications to the previously described ThermoMPNN framework¹⁰ (Fig. 1A).
93 ThermoMPNN is a structure-based protein stability model that extracts learned residue
94 embeddings from ProteinMPNN and passes these features through a lightweight prediction head
95 to obtain single mutant $\Delta\Delta G$ predictions. ProteinMPNN is a graph neural network trained to predict
96 protein sequences from the 3D structure of the protein¹⁵. Both models use message passing to
97 encode the local context surrounding the residue of interest, including the relative positions of
98 nearby residues. In this way, they use a combination of structure and sequence information to learn
99 what amino acids are likely to form favorable interactions if placed at a given position. In addition
100 to sequence and node embeddings from ProteinMPNN, we also extract directed edge features
101 representing the interaction between the mutated residue pair (Fig. 1B). We formulate our model
102 as a Siamese network by passing the concatenated per-mutation features through a shared
103 prediction head twice, once in each possible order. The raw predicted scores ($\Delta\Delta G_{AB}$ and $\Delta\Delta G_{BA}$)
104 are then symmetrized using a specialized loss function to enforce invariance to the mutation order
105 (details in Methods). We train ThermoMPNN-D on the double mutant subset of the Megascale
106 cDNA proteolysis dataset from Tsuboyama et al.¹⁶, which we call Megascale-D. Using this
107 scheme, ThermoMPNN-D obtains a high degree of order-invariance, with a Spearman correlation
108 coefficient (SCC) of 0.999 and average bias of 0.003 between $\Delta\Delta G_{AB}$ and $\Delta\Delta G_{BA}$ across the
109 Megascale-D test set.

110 Training ThermoMPNN-D on the Megascale-D dataset produced reasonable results on the
111 test split of the same dataset (SCC = 0.49 ± 0.01), but it struggled to generalize when tested on an

112 orthogonal test set from the literature, the Protherm double mutant dataset¹⁷, which we call
113 PTMUL-D (SCC = 0.35 ± 0.03) (Table 1, top section). After examining Megascale-D, we found
114 that, unlike its single mutant counterpart (Megascale-S), it is skewed in several ways. Most notably,
115 mutated residue pairs are typically close in 3D space, often in direct contact via side chain
116 interactions (Fig. 2A, blue bars), with a mean pairwise distance of 3.7Å. Wildtype residue pairs in
117 the dataset also tend to consist of large polar or aromatic groups engaged in strong couplings such
118 as hydrogen bonds and pi-cation interactions (Fig. 2B, blue bars). We hypothesized that training
119 on a dataset with these characteristics may lead to subpar generalizability. To address this issue,
120 we propose a new data augmentation trick which we call over-and-back data augmentation.

121

122 *Over-and-back data augmentation*

123 Our key observation is that every pair of single mutations in a protein are separated from each
124 other by two mutations. To construct an augmented data point (Fig. 1C), we select a single mutant
125 to serve as the wildtype state and invert its experimentally measured $\Delta\Delta G_{\text{single}}$ to represent the
126 reverse mutation. We then randomly sample a second single mutant within the same protein, but
127 at a different residue position, and add its experimentally measured $\Delta\Delta G_{\text{single}}$ to obtain our final
128 $\Delta\Delta G_{\text{double}}$. In this way, we can generate a much larger dataset which more evenly samples the
129 expected distribution in terms of pairwise distance and wildtype amino acid types (Figs. 2A and
130 2B, orange bars). In doing so, we hoped to enable our model to distinguish between distal, roughly
131 additive mutations and proximal, tightly coupled mutations. After retraining on the augmented
132 dataset, we observed significantly better results on PTMUL-D (SCC = 0.57 ± 0.02) at the cost of
133 a small drop in some Megascale-D metrics (Table 1, top panel). We noticed that this procedure
134 tends to generate a disproportionate fraction of stabilizing double mutants. Since most single

135 mutants are destabilizing, flipping the first $\Delta\Delta G_{\text{single}}$ tends to bias the resulting distribution toward
136 lower $\Delta\Delta G_{\text{double}}$ values (Fig. 2C, yellow peak). To partially correct for this effect, we implemented
137 a biased sampling procedure to shift the distribution closer to that of the non-augmented
138 Megascale-D dataset (Fig. 2C, orange peak). This adjustment slightly improved both root mean
139 squared error (RMSE) and correlation metrics across both datasets (Table 1, top panel).

140

141 *ThermoMPNN-D ablation study*

142 We next tested whether the Siamese aggregation scheme was necessary to achieve strong
143 performance (Table 1, middle panel). We found that this approach obtained better results on both
144 datasets when compared to previously proposed order-invariant aggregation functions such as
145 element-wise summation and averaging. We also experimented with modifying or removing other
146 components of our network (Table 1, bottom panel). We found that removing edge features slightly
147 degraded scores, but not as much as removing the Siamese aggregation. Additionally, we tested
148 fine-tuning ProteinMPNN by unfreezing the weights from the sequence recovery encoder/decoder,
149 which are kept fixed by default. Consistent with the original ThermoMPNN study, fine-tuning the
150 ProteinMPNN weights produced mixed results due to overfitting¹⁰. A small performance gain was
151 achieved by ensembling three independently trained models, a boost that we do not observe when
152 applied to single mutant ThermoMPNN. We suspect that this is enabled by the randomness
153 introduced by the data augmentation procedure. The final ensembled ThermoMPNN-D predictor
154 achieved SCC values of 0.54 and 0.59 on the Megascale-D and PTMUL-D test sets, respectively.

155

156 *Benchmarking ThermoMPNN-D against other double mutant models*

157 We then benchmarked ThermoMPNN-D against existing methods for double mutant $\Delta\Delta G$
158 prediction from the literature (Fig. 3). To do so, we performed 5-fold cross-validation across the
159 full Megascale-D dataset. We found that ThermoMPNN-D achieved state-of-the-art performance
160 on PTMUL-D, while only recent AlphaFold-based method Mutate Everything obtained
161 comparable performance on Megascale-D when evaluated on matching splits (SCC = 0.55). As a
162 baseline, we also included an additive ThermoMPNN prediction in which we simply added the
163 two predicted $\Delta\Delta G_{\text{single}}$ values for comparison to the epistasis-aware prediction of ThermoMPNN-
164 D. To our surprise, this method achieved even better results on Megascale-D (SCC = 0.59), along
165 with similar results on PTMUL-D, depending on the splits used. Intrigued by this finding, we
166 reevaluated each double mutant predictor from the literature by running a similar additive baseline
167 when available (Fig. 3A and 3B, green bars). We found that most methods provide little or no
168 improvement over their additive equivalent when utilized in epistatic mode. The only epistasis-
169 aware methods to provide better scores on both datasets were Rosetta and ESM-1v.

170 To further probe this phenomenon, we evaluated each predictor on Megascale-S for the
171 same set of proteins. We then plotted the single and double mutant error (RMSE) for each method
172 (Fig. 3C). All but two methods had lower error on single mutants, and they closely followed the
173 expected trajectory for the propagation of random additive errors. This indicates that the surveyed
174 methods generally fail to reduce the error on double mutants beyond what would be expected from
175 a purely additive model. The other two methods, FoldX and DDGun, instead followed the identity
176 line, with similar error on single and double mutants.

177 Since the Megascale dataset includes single and double mutant scans for the same proteins,
178 we can calculate the expected $\Delta\Delta G$ for a particular double mutant assuming an additive model
179 ($\Delta\Delta G_{\text{additive}}$). We plotted these values against the measured $\Delta\Delta G_{\text{double}}$ for the full Megascale-D

180 dataset (Fig. 2D). Notably, $\Delta\Delta G_{\text{double}}$ is highly correlated with $\Delta\Delta G_{\text{additive}}$ across the dataset (SCC
181 = 0.81), while the average observed epistatic coupling is -0.9 kcal/mol, indicating that $\Delta\Delta G_{\text{double}}$ is
182 typically less destabilizing than would be expected based on the observed $\Delta\Delta G_{\text{single}}$. Fitting a linear
183 regression to this dataset produces a y-intercept of 0.62 and a slope of 1.15, indicating that the
184 magnitude of epistatic effects increases with increasing $\Delta\Delta G_{\text{additive}}$.

185

186 *Deep mutational scan benchmark*

187 We next tested the same predictors on a collection of six deep mutational scans (DMS) gathered
188 from the literature (Table 2). Each DMS dataset consisted of at least 1,000 phenotypic
189 measurements for double mutants gathered in a single study (details in Table 3). Since these assays
190 each measure some proxy of protein fitness rather than stability, we anticipated lower correlations
191 with predicted $\Delta\Delta G$ than on the previous datasets. This was observed in most cases, and the best
192 methods across the full suite of assays were Mutate Everything (additive) and FoldX (epistatic),
193 with average SCC values of 0.40 and 0.39, respectively. Consistent with the prior results, most
194 methods show similar or worse performance in epistatic mode. Only FoldX produced equivalent
195 or better scores across all DMS assays.

196

197 *Stabilizing mutation detection*

198 We next evaluated stabilizing mutation predictions across the Megascale-D and PTMUL-D
199 datasets (Table 4). Stabilizing mutations are particularly difficult to predict, since the vast majority
200 of mutations are typically neutral or destabilizing compared to the wildtype. Indeed, less than 1%
201 of mutations in Megascale-D (n=1,254) fell under the commonly used threshold of $\Delta\Delta G \leq -0.5$
202 kcal/mol. Surprisingly, nearly every predictor showed improvement on both datasets when in

203 epistatic mode. While positive predictive value (PPV) showed mixed results in some cases, all
204 other metrics including Matthews Correlation Coefficient (MCC) consistently favored the epistatic
205 predictors. ThermoMPNN-D achieved the best prediction metrics on the Megascale-D and
206 PTMUL-D datasets, with an MCC of 0.19 and 0.38, respectively, compared to 0.17 and 0.28 for
207 additive ThermoMPNN. When evaluated on the cDNA2 test split of Megascale-D, Mutate
208 Everything (epistatic) outperforms ThermoMPNN-D (MCC = 0.27 vs 0.15), but the latter is more
209 effective on the PTMUL-D dataset when trained on the same splits (MCC = 0.38 for
210 ThermoMPNN-D vs 0.33 for Mutate Everything). We observe a significant discrepancy in
211 stabilizing mutation scores between PTMUL-D and Megascale-D, with nearly all methods
212 producing significantly better metrics on PTMUL-D in both additive and epistatic mode.
213

214 **Discussion**

215 This study was motivated by the hypothesis that a network designed to explicitly model double
216 point mutations could provide better $\Delta\Delta G$ predictions than a naïve model assuming additive
217 mutational effects. To test this hypothesis, we developed ThermoMPNN-D, which uses a Siamese
218 aggregation scheme and extensive data augmentation to leverage extensive mutagenesis data and
219 enforce helpful inductive biases such as the distance dependence of epistatic interactions and
220 mutation order invariance. Through rigorous benchmarking, we found our initial hypothesis was
221 not always correct, as ThermoMPNN-D and other double mutant predictors nearly all achieved
222 similar or worse overall results than their additive counterparts when evaluated by full-dataset
223 correlation coefficients. However, epistasis-aware predictors including ThermoMPNN-D enabled
224 improved prediction of stabilizing double mutations, which are critically important for protein
225 design applications.

226 Our study is one of the first to utilize the double mutant subset of the Megascale cDNA
227 proteolysis dataset recently published by Tsuboyama et al.¹⁸, which we call Megascale-D. As such,
228 it is important to note that models trained solely on Megascale-D proved unable to generalize to
229 unseen datasets. To address this issue, we introduce a novel data augmentation technique, over-
230 and-back augmentation, which may be considered as an extension of the recently introduced
231 thermodynamic permutation technique⁹ for sampling double mutations. The other extant study
232 utilizing the Megascale-D dataset also chose to expand their training dataset by pre-training on
233 Megascale-S¹⁹, although they did not evaluate a model trained only on Megascale-D. Taken
234 together, these findings raise the question: what constitutes a representative double mutant
235 landscape for modeling purposes? While exhaustive single mutant scans are now feasible for small
236 proteins, enumeration of double mutations remains challenging due to the exponential increase in

237 scale. With this in mind, we contend that data augmentation is an attractive strategy to expand the
238 pool of double mutant data to better capture the full mutational landscape. To enable further
239 development of data augmentation protocols, we make readily available our full dataset of 340,000
240 modeled mutant structures and Rosetta energies.

241 Most other protein stability models are limited to predicting single point mutations, while
242 even those with multiple mutation functionality have rarely been benchmarked against an
243 appropriate additive baseline. Still, a few previous studies provide evidence to corroborate our
244 findings. Consistent with our observations, Ouyang-Zhang et al. find that the epistatic version of
245 Mutate Everything behaves similarly to ThermoMPNN-D, in that its overall regression metrics are
246 similar or worse compared to its additive equivalent despite showing improved prediction of
247 stabilizing double mutations¹⁹. We also found that epistasis-aware models were often better
248 performing on certain datasets but worse on others. This is consistent with prior works which find
249 that epistatic terms derived from coevolutionary models are only beneficial for around 2/3 of tested
250 proteins^{20,21}, with factors such as MSA depth and assay design suggested as possible explanations.

251 We anticipated that predicting $\Delta\Delta G$ for double mutations would be more difficult than for
252 single mutations. This was generally observed, as top predictors including ThermoMPNN obtained
253 an SCC below 0.60 on Megascale-D, while the top reported score¹⁰ on Megascale-S is around
254 0.75. As expected, we also observe a lower success rate on stabilizing mutations, as ThermoMPNN
255 obtains a state-of-the-art PPV of 0.13 and 0.29 on different splits of Megascale-D compared to
256 0.45 on Megascale-S¹⁰. Double mutant data is less abundant than single mutant data, which makes
257 benchmarking more prone to random variance. To alleviate this issue, we employ DMS data to
258 supplement our stability datasets and cross-validate across all available data, which enabled
259 evaluation of >125,000 stability measurements and >74,000 DMS measurements gathered on

260 double mutations. Future work includes benchmarking and model development on higher-order
261 (3+) mutation datasets, which face even greater limitations in data availability and evaluation.

262 Epistasis is a complex phenomenon in which both global (per-protein) and local (per-
263 mutation) effects can influence variant fitness²², and their relative importance can vary by fitness
264 level and biological context²³. With this in mind, several avenues for future work may offer
265 potential for improvement. The pre-training schemes underpinning many recent models may be
266 redesigned to explicitly learn patterns of epistatic interaction rather than autoregressive or one-
267 shot decoding schemes. Model architecture may also be improved either by separating energetic
268 contributions from individual and pairwise residue terms, such as with a Potts model²⁴, or by
269 incorporating latent variables to represent global nonlinearities²⁵. Recent efforts to model protein
270 fitness with epistasis-aware neural networks^{26,27} may serve as a starting point for future protein
271 stability models. However, these methods tend to require parameterization with initial DMS data
272 for the target protein, so it remains to be seen how well they can generalize to novel proteins.

273

274 **Methods**

275 *ThermoMPNN-D architecture*

276 ThermoMPNN-D (Fig. 1A) was implemented as an extension of the ThermoMPNN framework¹⁰,
277 which uses sequence recovery model ProteinMPNN as a feature extractor¹⁵. All experiments used
278 the ProteinMPNN model trained with 0.2Å backbone noise, and ProteinMPNN weights were kept
279 frozen during training unless otherwise stated. To represent each mutation, we extracted the node
280 representation n_i for the mutated position from the molecular graph held in the last two decoder
281 layers of ProteinMPNN. We also retrieved the directed edge representation e_{ji} connecting from the
282 other mutated residue to the residue of interest (Fig. 1B). If no such edge existed (i.e., the mutations
283 are not within 48 nearest neighbors), a zero vector was substituted as the edge representation. We
284 subtracted the sequence embedding of the wildtype and mutant amino acids to obtain a sequence
285 representation s_i . The node, edge, and sequence representations were concatenated, and each
286 mutation vector was then passed through a shared MLP to aggregate and downsample to 128
287 dimensions. The mutation features were then concatenated in both AB and BA order, and each
288 permutation was passed through another shared MLP to produce raw predictions $\Delta\Delta G_{AB}$ and
289 $\Delta\Delta G_{BA}$, which were averaged to obtain a final $\Delta\Delta G$.

290

291 *ThermoMPNN-D training procedure*

292 ThermoMPNN-D includes 116,000 trainable parameters, which were trained for up to 100 epochs
293 using the Adam optimizer with an initial learning rate of 10^{-5} and a batch size of 256 mutations.
294 Dropout ($p=0.1$) and LayerNorm were used on all fully connected layers. Learning rate decay and
295 early stopping was conditioned on validation set mean squared error (MSE). Training used a

296 custom loss function inspired by antisymmetric single mutant predictor ACDC-NN²⁸ and applied
297 to the raw predictions $\Delta\Delta G_{AB}$ and $\Delta\Delta G_{BA}$:

$$298 \quad loss = MSE(\Delta\Delta G_{true}, \Delta\Delta G_{avg}) + \langle \Delta\Delta G_{sym} \rangle$$

$$299 \quad \Delta\Delta G_{avg} = \frac{\Delta\Delta G_{AB} + \Delta\Delta G_{BA}}{2}$$

$$300 \quad \langle \Delta\Delta G_{sym} \rangle = \left\langle \frac{|\Delta\Delta G_{AB} - \Delta\Delta G_{BA}|}{2} \right\rangle$$

301 A non-Siamese model was built to test other aggregators (Table 1, middle panel). This model used
302 the same featurization scheme, but after downsampling, mutation embeddings were aggregated
303 instead of concatenated and passed once through the final MLP. Fine-tuning ProteinMPNN was
304 implemented by unfreezing all layers with a separate learning rate, which was selected via
305 parameter sweep (10^{-6} gave the best results). Ensembling was implemented by averaging the
306 predicted $\Delta\Delta G$ from three independently trained models with different random seeds for training
307 and data augmentation.

308

309 *Over-and-back data augmentation*

310 For each single mutant in the Megascale training set, the modeled mutant structure was obtained
311 using Rosetta¹¹. The second mutation was sampled stochastically from all possible single
312 mutations that a) shared the same PDB ID and b) did not share the same amino acid position. To
313 bias sampling toward more destabilizing $\Delta\Delta G$ values, the $\Delta\Delta G_{single}$ values for the whole dataset
314 were used to obtain a weighted sampling probability (P) as follows:

$$315 \quad y = -1 * \Delta\Delta G_{true}$$

$$316 \quad P = [y - \min(y)]^3$$

317 This distribution was normalized for each individual mutation. Augmented datasets were sampled
318 once at the beginning of training and randomly shuffled after every epoch.

319

320 *Dataset splits and curation*

321 For the ThermoMPNN-D ablation study, we obtained the Megascale dataset reported in
322 Tsuboyama et al.¹⁸ from its Zenodo repository¹⁶, following the splitting procedure previously
323 described for ThermoMPNN¹⁰, with the following modifications. We removed any homologues
324 (>25% sequence similarity) to proteins in the PTMUL dataset. Second, we trained on double
325 mutants with defined ddG_{ML} values. After removing duplicate data points, we obtained a
326 training/validation/test split of 85,253/10,282/18,574 mutations across 90/17/20 proteins.

327 For the double mutant model benchmarks, we used the full Megascale dataset and
328 evaluated ThermoMPNN using 5-fold cross-validation split by sequence similarity, as previously
329 described. To compare additive and epistatic models, we matched single and double mutant
330 measurements and dropped any double mutants without valid single mutant data, resulting in
331 127,476 double mutations across 153 proteins. The Protherm multiple mutation (PTMUL) dataset
332 introduced in the DDGun paper¹⁷ and re-curated for Mutate Everything¹⁹ was used after dropping
333 higher-order (3+) mutation measurements, resulting in 536 mutations across 83 proteins (PTMUL-
334 D). Since Mutate Everything was trained on different splits of the Megascale dataset, we retrained
335 and reevaluated ThermoMPNN using their training/test splits, which they denote “cDNA2”,
336 resulting in a test set of 22,913 mutations across 18 proteins. For the single vs double mutant error
337 calculation, we used the full single mutant Megascale dataset (Megascale-S), which contained
338 271,231 mutations across 298 proteins.

339 We curated deep mutational scanning (DMS) datasets from the ProteinGym benchmark²⁹.
340 We selected DMS datasets with >1000 double mutations and endpoints that might serve as
341 reasonable proxies for thermodynamic stability. From this pool, we eliminated assays overlapping
342 with the Megascale dataset and those without a high-confidence AlphaFold model or crystal
343 structure. We were then left with six assays, which are summarized in Table 3.

344

345 *Literature model benchmarking*

346 For the Rosetta benchmark, we adapted a previously published *RosettaScripts* point mutation
347 protocol¹¹ for use on double mutations by applying constraints to all residues nearby to either
348 residue. To convert REU into approximate kcal/mol units, we divided all energy values by 2.9, as
349 recommended for the ref2015 score function³⁰. FoldX was downloaded under an academic license
350 (<https://foldxsuite.crg.eu>), and predictions were obtained by running RepairPDB on all input
351 structures, followed by PositionScan for single mutants or additive predictions and BuildModel
352 for epistatic predictions³¹. MAESTRO³² was downloaded from its website
353 (<https://pbwww.services.came.sbg.ac.at>), while DDGun/DDGun3D¹⁷
354 (<https://github.com/biofold/ddgun>), ESM-1v³³ (<https://github.com/facebookresearch/esm>),
355 ProteinMPNN¹⁵ (<https://github.com/dauparas/ProteinMPNN>), and Mutate Everything¹⁹
356 (<https://github.com/jozhang97/MutateEverything>) were obtained from their respective GitHub
357 repositories.

358 ProteinMPNN zero-shot predictions were obtained by masking out the mutated residue(s)
359 and calculating the difference in negative log-likelihood between the mutant and wildtype residues.
360 For the ESM zero-shot predictions, we used an ensemble of five ESM-1v (650M, UR90S) models
361 with the masked-marginals scoring method, as recommended³³. To obtain epistatic predictions for

362 ProteinMPNN and ESM-1v, both mutated residues were masked prior to inference, while the
363 additive predictions masked each residue individually.

364

365 *Theoretical error calculation*

366 We calculated the theoretical error for double mutant predictions as follows:

367
$$\sigma_{AB} = \sqrt{\sigma_A^2 + \sigma_B^2}$$

368 Where σ_A and σ_B are the single mutant prediction errors (in RMSE) for mutation A and B, and σ_{AB}
369 is the theoretical error for double mutants. Note that this model assumes that single mutant errors
370 are randomly distributed and uncorrelated.

371

372 *Stabilizing mutation metrics*

373 To evaluate stabilizing mutation predictions (Table 4), we primarily use the Matthews correlation
374 coefficient (MCC), which is widely accepted as a robust holistic measure of classifier accuracy on
375 unbalanced datasets³⁴. Following the convention from Ouyang-Zhang et al.¹⁹, we calculate MCC
376 across the full dataset using a threshold of 0 kcal/mol. For the remaining metrics, we use the
377 definition that mutations with $\Delta\Delta G \leq -0.5$ kcal/mol are stabilizing. This resulted in 1254, 111, and
378 198 stabilizing mutations for the Megascale-D, PTMUL-D, and cDNA2 test datasets, respectively.

379 We calculate the positive predictive value (PPV) across each full dataset, while detection
380 precision (DetPr) and normalized discounted cumulative gain (nDCG) are calculated separately
381 for each protein and averaged. To calculate these last two metrics, the mutations for a given protein
382 are sorted by predicted $\Delta\Delta G$, and the top K mutations are selected (K=30 in this study). The DetPr
383 represents the fraction of top-30 mutations that are measured to be truly stabilizing, while nDCG
384 is a more complicated measure of how highly the model ranks the best 30 mutations.

385

386 **Code Availability**

387 ThermoMPNN-D trained model weights and code are available at [https://github.com/Kuhlman-](https://github.com/Kuhlman-Lab/ThermoMPNN-D)
388 [Lab/ThermoMPNN-D](https://github.com/Kuhlman-Lab/ThermoMPNN-D).

389

390 **Data Availability**

391 The full Megascale dataset can be obtained from its Zenodo repository¹⁶, while the full
392 ProteinGym datasets are available at <https://proteingym.org>²⁹ and the full PTMUL dataset is
393 available at <https://github.com/jozhang97/MutateEverything>. The curated Megascale, PTMUL-D,
394 and DMS double mutant datasets and splits used in this study are available on Zenodo at
395 <https://doi.org/10.5281/zenodo.13345274>. Modeled single mutant structures and energies obtained
396 using Rosetta for the full Megascale dataset are available in the same repository.

397

398 **Supplementary material description:**

399 N/A

400

401 **Conflict of interest statement:**

402 The authors have no relevant conflicts of interest to declare.

403

404 **Acknowledgements:**

405 This work was supported by NIH grant R35GM131923 (B.K.) and NSF fellowship DGE-2040435
406 (H.D.). H.D. acknowledges support by a Pre-doctoral Fellowship from the American Foundation
407 for Pharmaceutical Education. This work utilized the resources of the UNC Longleaf high-

408 performance computing cluster. The authors would like to thank Dr. Pranam Chatterjee for his
409 advice regarding protein language models and Jeffrey O. Zhang for his assistance with the Mutate
410 Everything platform.

411 **References:**

- 412 1. Zheng J, Guo N, Wagner A (2020) Selection enhances protein evolvability by increasing
413 mutational robustness and foldability. *Science* 370.
- 414 2. Høie MH, Cagiada M, Beck Frederiksen AH, Stein A, Lindorff-Larsen K (2022) Predicting and
415 interpreting large-scale mutagenesis data using analyses of protein stability and conservation. *Cell*
416 *Rep.* 38:110207.
- 417 3. Hartl FU (2017) Protein Misfolding Diseases. *Annu. Rev. Biochem.* 86:21–26.
- 418 4. Sawaya MR, Hughes MP, Rodriguez JA, Riek R, Eisenberg DS (2021) The expanding amyloid
419 family: Structure, stability, function, and pathogenesis. *Cell* 184:4857–4873.
- 420 5. Narayanan H, Dingfelder F, Butté A, Lorenzen N, Sokolov M, Arosio P (2021) Machine
421 learning for biologics: opportunities for protein engineering, developability, and formulation.
422 *Trends Pharmacol. Sci.* 42:151–165.
- 423 6. Zhu Z, Song H, Wang Y, Zhang Y-HP (2022) Protein engineering for electrochemical
424 biosensors. *Curr. Opin. Biotechnol.* 76:102751.
- 425 7. Notin P, Rollins N, Gal Y, Sander C, Marks D (2024) Machine learning for functional protein
426 design. *Nat. Biotechnol.* 42:216–228.
- 427 8. Chen Y, Xu Y, Liu D, Xing Y, Gong H (2024) SPIRED-Fitness: an end-to-end framework for
428 the prediction of protein structure and fitness from single sequence. *BioRxiv*.
- 429 9. Diaz DJ, Gong C, Ouyang-Zhang J, Loy JM, Wells J, Yang D, Ellington AD, Dimakis A,
430 Klivans AR (2023) Stability Oracle: A Structure-Based Graph-Transformer for Identifying
431 Stabilizing Mutations. *BioRxiv*.
- 432 10. Dieckhaus H, Brocidiaco M, Randolph NZ, Kuhlman B (2024) Transfer learning to leverage
433 larger datasets for improved prediction of protein stability changes. *Proc Natl Acad Sci USA*
434 121:e2314853121.
- 435 11. Thieker DF, Maguire JB, Kudlacek ST, Leaver-Fay A, Lyskov S, Kuhlman B (2022)
436 Stabilizing proteins, simplified: A Rosetta-based webtool for predicting favorable mutations.
437 *Protein Sci.* 31:e4428.
- 438 12. Saito Y, Koya J, Kataoka K (2021) Multiple mutations within individual oncogenes. *Cancer*
439 *Sci.* 112:483–489.
- 440 13. Lou H, Chen M, Black SS, Bushell SR, Ceccarelli M, Mach T, Beis K, Low AS, Bamford VA,
441 Booth IR, et al. (2011) Altered antibiotic transport in OmpC mutants isolated from a series of
442 clinical strains of multi-drug resistant *E. coli*. *PLoS ONE* 6:e25825.

- 443 14. Kumar R, Srivastava Y, Muthuramalingam P, Singh SK, Verma G, Tiwari S, Tandel N, Beura
444 SK, Panigrahi AR, Maji S, et al. (2023) Understanding Mutations in Human SARS-CoV-2 Spike
445 Glycoprotein: A Systematic Review & Meta-Analysis. *Viruses* 15.
- 446 15. Dauparas J, Anishchenko I, Bennett N, Bai H, Ragotte RJ, Milles LF, Wicky BIM, Courbet A,
447 de Haas RJ, Bethel N, et al. (2022) Robust deep learning-based protein sequence design using
448 ProteinMPNN. *Science* 378:49–56.
- 449 16. Tsuboyama K, Dauparas J, Chen J, Laine E, Behbahani YM, Weinstein JJ, Mangan NM,
450 Ovchinnikov S, Rocklin GJ (2023) Mega-scale experimental analysis of protein folding stability
451 in biology and protein design. Zenodo.
- 452 17. Montanucci L, Capriotti E, Frank Y, Ben-Tal N, Fariselli P (2019) DDGun: an untrained
453 method for the prediction of protein stability changes upon single and multiple point variations.
454 *BMC Bioinformatics* 20:335.
- 455 18. Tsuboyama K, Dauparas J, Chen J, Laine E, Mohseni Behbahani Y, Weinstein JJ, Mangan
456 NM, Ovchinnikov S, Rocklin GJ (2023) Mega-scale experimental analysis of protein folding
457 stability in biology and design. *Nature* 620:434–444.
- 458 19. Ouyang-Zhang J, Diaz D, Klivans A, Kraehenbuehl P (2023) Predicting a Protein’s Stability
459 under a Million Mutations. *Advances in Neural Information Processing Systems*.
- 460 20. Laine E, Karami Y, Carbone A (2019) GEMME: A simple and fast global epistatic model
461 predicting mutational effects. *Mol. Biol. Evol.* 36:2604–2619.
- 462 21. Hopf TA, Ingraham JB, Poelwijk FJ, Schärfe CPI, Springer M, Sander C, Marks DS (2017)
463 Mutation effects predicted from sequence co-variation. *Nat. Biotechnol.* 35:128–135.
- 464 22. Johnson MS, Reddy G, Desai MM (2023) Epistasis and evolution: recent advances and an
465 outlook for prediction. *BMC Biol.* 21:120.
- 466 23. Otwinowski J, McCandlish DM, Plotkin JB (2018) Inferring the shape of global epistasis. *Proc*
467 *Natl Acad Sci USA* 115:E7550–E7558.
- 468 24. Li AJ, Lu M, Desta I, Sundar V, Grigoryan G, Keating AE (2023) Neural network-derived
469 Potts models for structure-based protein design using backbone atomic coordinates and tertiary
470 motifs. *Protein Sci.* 32:e4554.
- 471 25. Ding D, Shaw AY, Sinai S, Rollins N, Prywes N, Savage DF, Laub MT, Marks DS (2024)
472 Protein design using structure-based residue preferences. *Nat. Commun.* 15:1639.
- 473 26. Luo Y, Jiang G, Yu T, Liu Y, Vo L, Ding H, Su Y, Qian WW, Zhao H, Peng J (2021) ECNet
474 is an evolutionary context-integrated deep learning framework for protein engineering. *Nat.*
475 *Commun.* 12:5743.

- 476 27. Aghazadeh A, Nisonoff H, Ocal O, Brookes DH, Huang Y, Koyluoglu OO, Listgarten J,
477 Ramchandran K (2021) Epistatic Net allows the sparse spectral regularization of deep neural
478 networks for inferring fitness functions. *Nat. Commun.* 12:5225.
- 479 28. Benevenuta S, Pancotti C, Fariselli P, Birolo G, Sanavia T (2021) An antisymmetric neural
480 network to predict free energy changes in protein variants. *J. Phys. D Appl. Phys.* 54:245403.
- 481 29. Notin P, Kollasch AW, Ritter D, van Niekerk L, Paul S, Spinner H, Rollins N, Shaw A,
482 Weitzman R, Frazer J, et al. (2023) ProteinGym: Large-Scale Benchmarks for Protein Design and
483 Fitness Prediction. *BioRxiv*.
- 484 30. Park H, Bradley P, Greisen P, Liu Y, Mulligan VK, Kim DE, Baker D, DiMaio F (2016)
485 Simultaneous Optimization of Biomolecular Energy Functions on Features from Small Molecules
486 and Macromolecules. *J. Chem. Theory Comput.* 12:6201–6212.
- 487 31. Schymkowitz J, Borg J, Stricher F, Nys R, Rousseau F, Serrano L (2005) The FoldX web
488 server: an online force field. *Nucleic Acids Res.* 33:W382-8.
- 489 32. Laimer J, Hofer H, Fritz M, Wegenkittl S, Lackner P (2015) MAESTRO--multi agent stability
490 prediction upon point mutations. *BMC Bioinformatics* 16:116.
- 491 33. Meier J, Rao R, Verkuil R, Liu J, Sercu T, Rives A (2021) Language models enable zero-shot
492 prediction of the effects of mutations on protein function. *BioRxiv*.
- 493 34. Gösgens M, Zhiyanov A, Tikhonov A, Prokhorenkova L (2021) Good Classification Measures
494 and How to Find Them. *Advances in Neural Information Processing Systems*.
- 495 35. Sarkisyan KS, Bolotin DA, Meer MV, Usmanova DR, Mishin AS, Sharonov GV, Ivankov DN,
496 Bozhanova NG, Baranov MS, Soylemez O, et al. (2016) Local fitness landscape of the green
497 fluorescent protein. *Nature* 533:397–401.
- 498 36. Gonzalez Somermeyer L, Fleiss A, Mishin AS, Bozhanova NG, Igoikina AA, Meiler J, Alaball
499 Pujol M-E, Putintseva EV, Sarkisyan KS, Kondrashov FA (2022) Heterogeneity of the GFP fitness
500 landscape and data-driven protein design. *eLife* 11.
- 501 37. Pokusaeva VO, Usmanova DR, Putintseva EV, Espinar L, Sarkisyan KS, Mishin AS,
502 Bogatyreva NS, Ivankov DN, Akopyan AV, Avvakumov SY, et al. (2019) An experimental assay
503 of the interactions of amino acids from orthologous sequences shaping a complex fitness
504 landscape. *PLoS Genet.* 15:e1008079.
- 505 38. Weng C, Faure AJ, Escobedo A, Lehner B (2024) The energetic and allosteric landscape for
506 KRAS inhibition. *Nature* 626:643–652.
- 507
- 508

509 **Tables:**

510 **Table 1:** ThermoMPNN-D ablation study results. D.A. stands for data augmentation.

Trial	Megascale-D			PTMUL-D		
	PCC	SCC	RMSE	PCC	SCC	RMSE
No D.A.	0.54 ± 0.02	0.49 ± 0.01	0.96 ± 0.01	0.36 ± 0.04	0.35 ± 0.03	2.11 ± 0.02
Naïve D.A.	0.50 ± 0.02	0.51 ± 0.02	1.19 ± 0.03	0.55 ± 0.03	0.57 ± 0.02	2.06 ± 0.04
Biased D.A.	0.52 ± 0.02	0.53 ± 0.02	1.09 ± 0.02	0.55 ± 0.02	0.57 ± 0.02	1.96 ± 0.03
Siamese	0.52 ± 0.02	0.53 ± 0.02	1.09 ± 0.02	0.55 ± 0.02	0.57 ± 0.02	1.96 ± 0.03
Max	0.50 ± 0.01	0.52 ± 0.01	1.15 ± 0.02	0.50 ± 0.02	0.54 ± 0.01	2.05 ± 0.02
Mean	0.43 ± 0.01	0.42 ± 0.01	1.19 ± 0.02	0.50 ± 0.01	0.52 ± 0.01	2.02 ± 0.01
Sum	0.45 ± 0.01	0.43 ± 0.01	1.19 ± 0.03	0.50 ± 0.01	0.53 ± 0.01	2.02 ± 0.02
Product	0.46 ± 0.04	0.47 ± 0.03	1.27 ± 0.02	0.49 ± 0.03	0.52 ± 0.02	2.07 ± 0.03
Baseline	0.52 ± 0.02	0.53 ± 0.02	1.09 ± 0.02	0.55 ± 0.02	0.57 ± 0.02	1.96 ± 0.03
- Edges	0.49 ± 0.01	0.51 ± 0.01	1.13 ± 0.01	0.52 ± 0.01	0.56 ± 0.02	2.00 ± 0.01
+ Fine-tune	0.47 ± 0.02	0.48 ± 0.02	1.15 ± 0.01	0.55 ± 0.01	0.59 ± 0.01	1.96 ± 0.03
+ Ensemble	0.54	0.55	1.07	0.57	0.59	1.95

511 All statistics are reported as mean ± s.d. of triplicate runs, except for the ensemble.

512 **Table 2:** Deep mutational scan benchmark results for selected double mutant prediction methods
 513 (additive/epistatic models). The score of the best method on each assay is bolded.

Model	Spearman Correlation Coefficient						Mean
	avGFP	cgreGFP	ppluGFP2	amacGFP	His3	KRas	
Rosetta ¹¹	0.41/0.42	0.37/0.34	0.29/0.29	0.21/0.21	0.26/0.26	0.37/0.35	0.32/0.31
FoldX ³¹	0.46/0.47	0.52/0.52	0.39/0.39	0.38/0.38	0.20/0.26	0.34/0.34	0.38/0.39
DDGun ¹⁷	0.13/--	0.32/--	0.17/--	0.14/--	0.14/--	0.21/--	0.19/--
DDGun3D ¹⁷	0.27/--	0.31/--	0.18/--	0.16/--	0.11/--	0.24/--	0.21/--
MAESTRO ³²	0.26/0.22	0.23/0.15	0.14/0.08	0.11/0.07	0.17/0.13	0.25/0.26	0.19/0.15
ESM-1v ³³	0.00/0.01	0.01/0.02	-0.01/0.02	-0.01/0.01	0.14/0.21	0.19/0.20	0.05/0.08
ProteinMPNN ¹⁵	0.35/0.36	0.23/0.26	0.12/0.11	0.13/0.12	0.18/0.15	0.36/0.37	0.23/0.22
ThermoMPNN	0.46/0.40	0.40/0.24	0.21/0.03	0.26/0.16	0.28/0.24	0.37/0.31	0.33/0.23
ThermoMPNN*	0.48/0.44	0.40/0.22	0.21/0.06	0.28/0.18	0.29/0.24	0.39/0.31	0.34/0.24
Mutate	0.53/0.49	0.50/0.43	0.37/0.30	0.32/0.27	0.27/0.27	0.40/0.36	0.40/0.35
Everything ¹⁹							

514 * Retrained on cDNA training splits from Ouyang-Zhang et al.¹⁹

515

516 **Table 3:** Summary of curated deep mutational scan assays of double mutants.

Assay name and source	Abbreviation	Mutations	Phenotype
GFP_AEQVI ³⁵	avGFP	12,777	Fluorescence
D7PM05_CLYGR ³⁶	cgreGFP	10,148	Fluorescence
Q6WV13_9MAXI ³⁶	ppluGFP2	15,992	Fluorescence
Q8WTC7_9CNID ³⁶	amacGFP	11,260	Fluorescence
HIS7_YEAST ³⁷	His3	1,475	Enzyme activity
RASK_HUMAN ³⁸	KRas	22,946	Expression

517

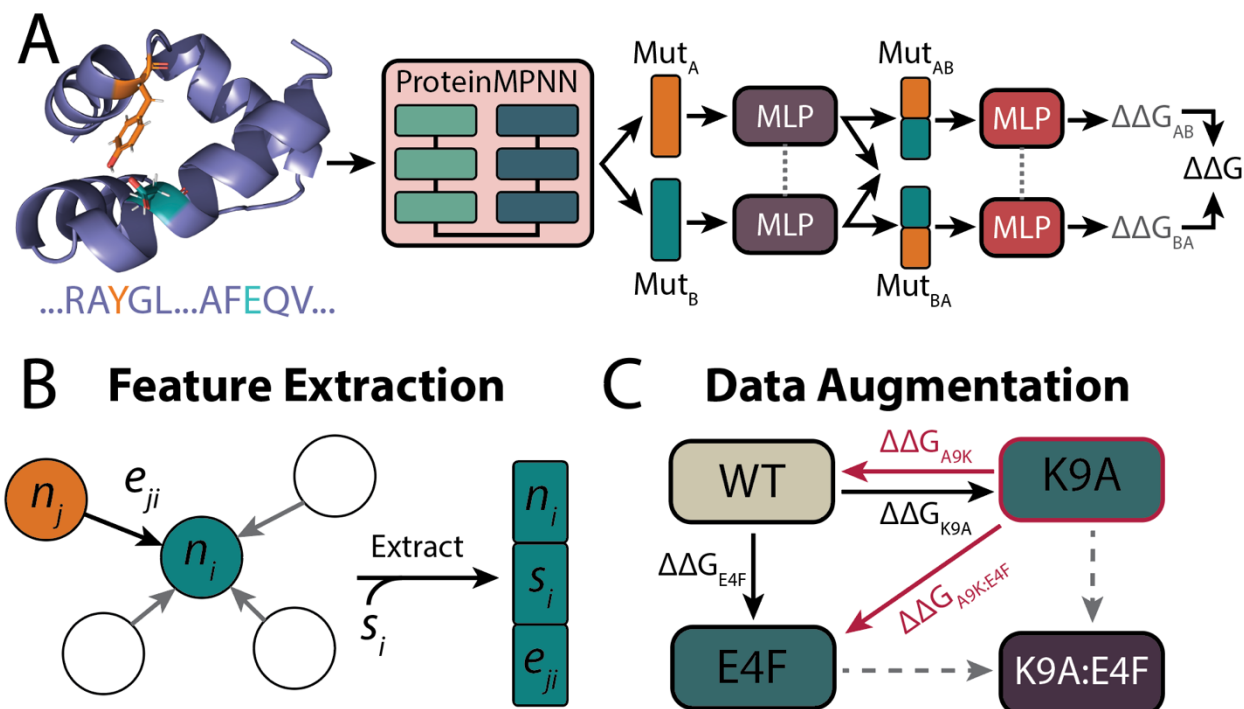
518 **Table 4:** Stabilizing mutation detection metrics for selected prediction methods (additive/epistatic
 519 models). The score of the best method on each metric is bolded.

Model	Megascale-D (n=1,254)				PTMUL-D (n=111)	
	MCC	PPV	DetPr ₃₀	nDCG ₃₀	MCC	PPV
Rosetta ¹¹	0.11/0.15	0.05/0.09	0.07/0.12	0.15/0.20	0.29/0.29	0.48/0.42
FoldX ³¹	0.13/0.14	0.04/0.04	0.07/0.08	0.16/0.16	0.22/0.24	0.38/0.36
DDGun ¹⁷	0.12/--	0.04/--	0.10/--	0.18/--	0.22/--	0.50/--
DDGun3D ¹⁷	0.13/--	0.05/--	0.08/--	0.17/--	0.17/--	0.47/--
MAESTRO ³²	0.15/0.14	0.04/0.03	0.09/0.09	0.13/0.17	--/--	--/--
ESM-1v ³³	0.02/0.03	0.01/0.02	0.03/0.05	0.05/0.12	0.07/0.09	0.30/0.31
ProteinMPNN ¹⁵	0.07/0.10	0.06/0.05	0.07/0.09	0.17/0.19	0.30/0.33	0.51/0.49
ThermoMPNN	0.17/ 0.19	0.13/0.13	0.20/ 0.22	0.31/ 0.35	0.29/ 0.37	0.49/ 0.57
	cDNA2 test (n=198)				PTMUL-D (n=111)	
ThermoMPNN*	0.10/0.15	0.29/0.20	0.10/0.17	0.11/0.22	0.34/ 0.38	0.58/0.54
Mutate Everything ¹⁹	0.26/ 0.27	0.12/0.11	0.24/ 0.30	0.35/ 0.43	0.33/0.33	0.46/0.44

520 * Re-trained on cDNA training splits from Ouyang-Zhang et al.¹⁹

521

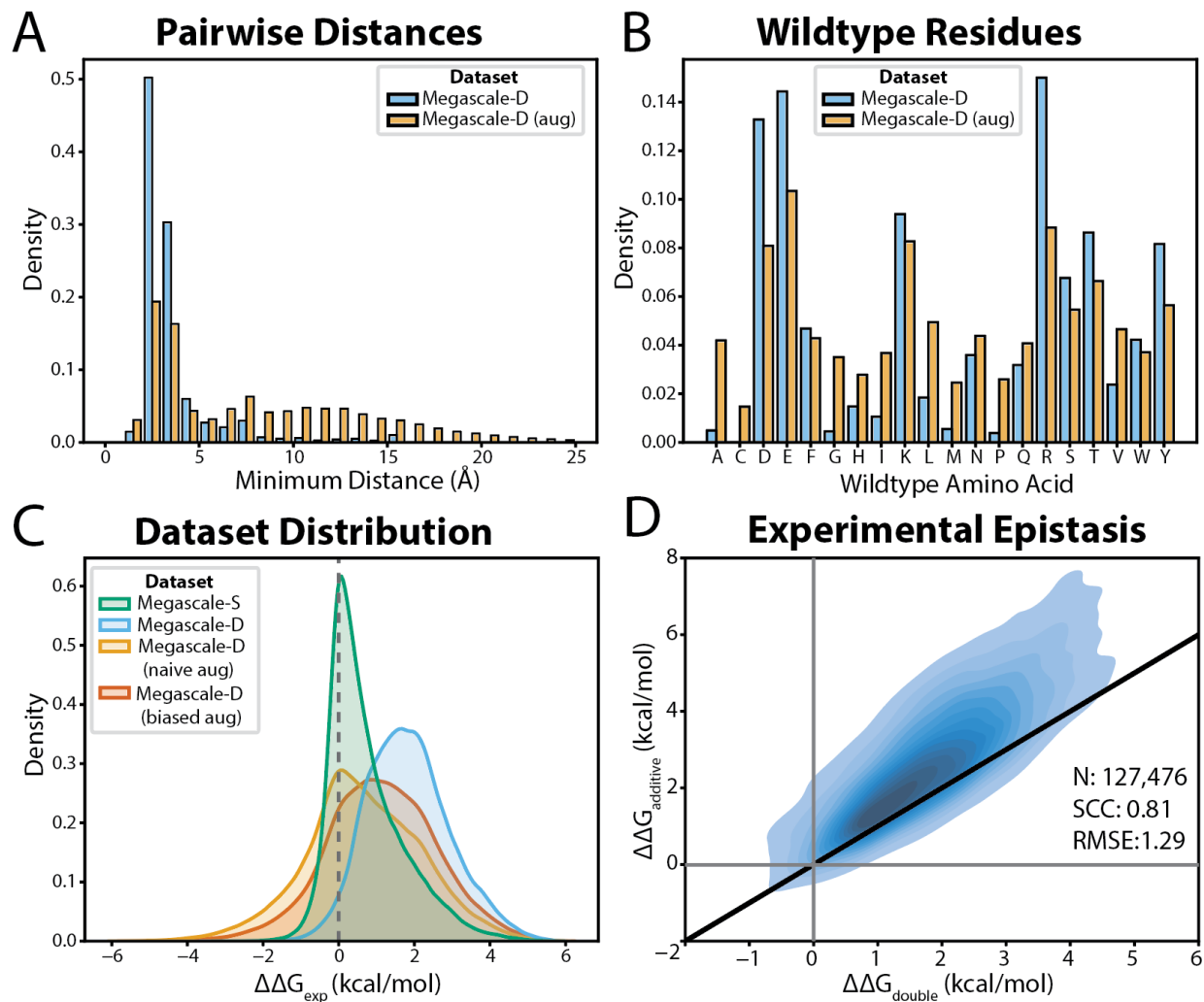
522 **Figures:**



523

524 **Figure 1: The ThermoMPNN-D modeling framework.** A) Schematic of ThermoMPNN-D, a
525 Siamese neural network for predicting double mutant stability changes. Dashed grey lines indicate
526 shared weights. B) Example feature extraction step for hypothetical mutation i , in which the node
527 (n_i), sequence (s_i), and edge (e_{ji}) embeddings are extracted from the protein graph. C)
528 Thermodynamic cycle demonstrating the principle of over-and-back data augmentation. Black
529 arrows denote mutations with a defined $\Delta\Delta G$ in the original dataset, dashed grey arrows indicate
530 mutations missing data, and red arrows indicate mutations defined only via augmentation. The
531 augmented wildtype state is outlined in red.

532



533

534 **Figure 2: Megascale double mutant (Megascale-D) dataset analysis and augmentation. A)**

535 Frequency of mutations stratified by minimum pairwise interatomic distance between mutated

536 residues and B) frequency of wildtype amino acids in the original and augmented Megascale-D.

537 C) Kernel density estimate distributions of Megascale dataset $\Delta\Delta G$ values with and without

538 augmentation. Dashed grey line indicates a theoretical neutral mutation. More positive $\Delta\Delta G$ values

539 indicate more destabilizing mutations. D) Kernel density estimate plot of Megascale-D comparing

540 measured double mutant $\Delta\Delta G$ to the corresponding additive $\Delta\Delta G$ obtained from the sum of the

541 two constituent single mutants. The identity line is shown in black.

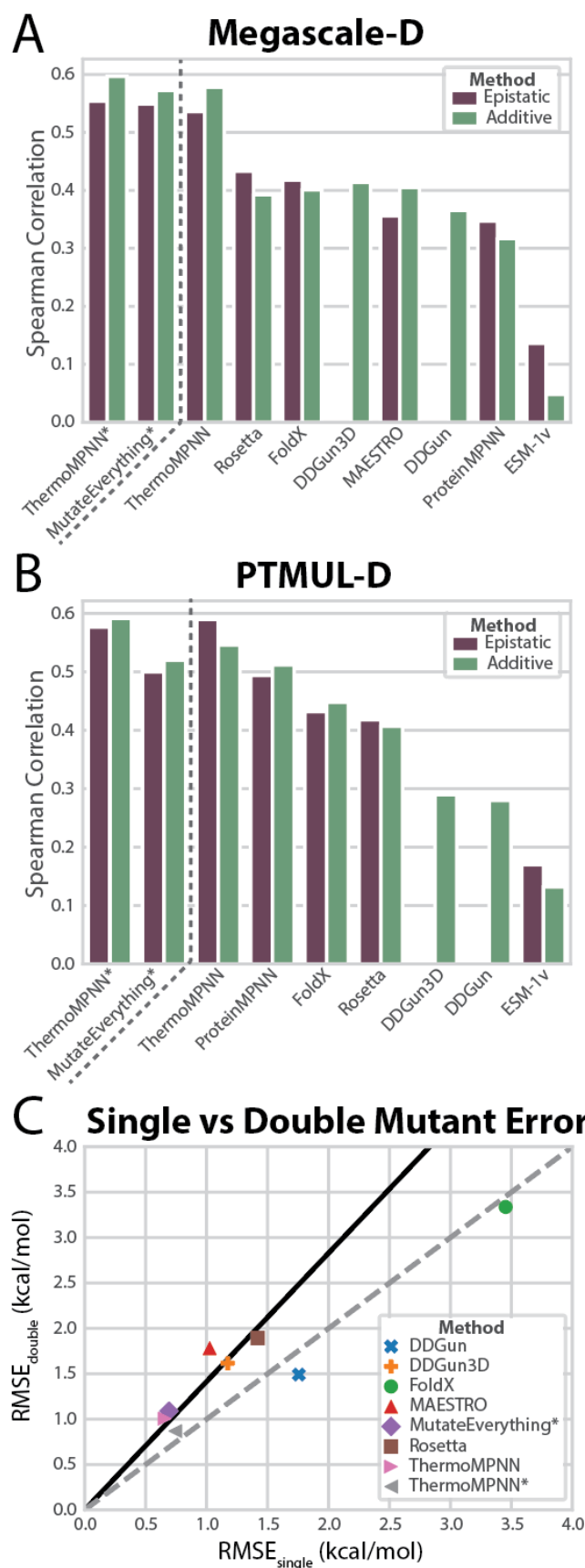


Figure 3: Comparison of ThermoMPNN and selected prior methods for modeling double mutants. A-B) Spearman correlation of selected additive and epistatic methods on A) the Megascale double mutant dataset (N=127,476) and B) the PTMUL double mutant dataset (N=536). Methods marked with asterisks were retrained and evaluated using different Megascale dataset splits. C) Root mean squared error (RMSE) of selected methods on the Megascale single mutant (x-axis) and double mutant (y-axis) datasets. The identity line is shown in dashed grey, and the theoretical error for a method following naïve additive error propagation behavior is shown in solid black.