
Designing active and thermostable enzymes with sequence-only predictive models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Data-driven models of protein fitness can be useful in designing novel proteins with
2 improved properties, but many questions remain regarding how and in what settings
3 they should be used. Here, we ask: How can we use predictive models of protein
4 fitness, whose predictions we might not always trust, to design protein sequences
5 enhanced for multiple fitness functions? We propose a general approach for doing
6 so, and apply it to design novel variants of eight different acylphosphatase and
7 lysozyme wild types, intended to be more thermostable and at least as catalytically
8 active as the wild types. Our method does not require a structure, experimental
9 measurements of activity, curation of homologous sequences, or family-specific
10 thermostability data. Experimental characterizations of our designed sequences,
11 as well as sequences designed by PROSS, a competitive baseline method for
12 improving protein thermostability, are currently underway and forthcoming.

13 1 Introduction

14 The design of proteins with improved fitness has historically been accomplished in the wet lab with
15 resource-intensive directed evolution campaigns, or *in silico* with compute-intensive biophysics
16 simulations. Recent years, however, have seen increasing interest in augmenting or replacing
17 such methods with data-driven approaches. One paradigm is to build a data-driven model of the
18 fitness function of interest—for example, fluorescence or enzyme activity (1)—then deploy some
19 optimization algorithm to identify novel sequences predicted by the model to exhibit higher fitness
20 than known protein sequences (1; 2; 3; 4).

21 In this work, we aim to make progress on two open challenges in this paradigm. First, how can we
22 design sequences that are enhanced for *multiple* fitness functions? Second, how can we leverage
23 predictive models for this task, while accounting for the fact that they are not equally trustworthy
24 everywhere in sequence space (5)? We propose one general approach for doing so, and apply it to
25 design novel enzyme variants intended to be more thermostable and at least as catalytically active as
26 wild type (WT) enzymes on their native reactions, for eight WTs from two different enzyme families:
27 acylphosphatases (ACYPs) and lysozymes.

28 Engineering thermostable enzymes is a longstanding problem of broad interest, as it enables their
29 deployment in the industrial conditions required for the production of food, biofuels, and pharma-
30 ceuticals, wastewater treatment, and a host of other industrial applications (6). However, it remains
31 challenging, in part due to inherent trade-offs between enzyme thermostability and activity (7; 8).
32 Our *in silico* results are promising: the melting temperatures of our designed variants are predicted to
33 surpass those of both the respective WTs and variants designed by PROSS, a competitive baseline
34 method for enhancing thermostability (9; 10), and for six out of the eight WTs, our designed variants
35 are also predicted to be at least as active as the WT. More importantly, experimental characterization
36 of the thermostabilities and activities of our designed sequences and those designed by PROSS are cur-

37 rently underway. We look forward to sharing the forthcoming data, which will reveal whether—and
 38 more importantly, why—our proposed approach did or did not succeed.

39 Another goal of this work was to investigate the use of likelihoods of models such as PROGEN2
 40 (11)—high-capacity models of the distribution of natural protein sequences, which we henceforth
 41 refer to as *natural protein sequence density models*—as a proxy for the catalytic activity of enzymes
 42 on their native reactions, in order to design active enzyme sequences. Likelihoods of such models
 43 have been successfully used to design catalytically active lysozymes (12) and high-affinity antibodies
 44 without the need for fitness-labeled data (13). We build upon these precedents by pursuing the
 45 enhancement of multiple fitness functions—catalytic activity and thermostability—for two very
 46 different enzyme families.

47 Finally, our predictive model for thermostability was trained on the “meltome” dataset, which contains
 48 the melting temperatures of more than 34k unique protein sequences from the proteomes of thirteen
 49 species (14), from which we excluded any sequences from our two families of interest (nine ACYP
 50 sequences and no lysozyme sequences). If our designed sequences prove to be both catalytically
 51 active and thermostable in the wet lab, commensurate with predictions, it means our method can
 52 be used to design active, thermostable enzymes without a structure, experimental measurements of
 53 activity, curation of homologous sequences, or family-specific thermostability data.

54 2 Proposed computational approach

55 Let \mathcal{X} denote the set of all amino acid sequences of length L . We consider the general goal of
 56 generating sequences, $x \in \mathcal{X}$, which are as diverse as possible, while exhibiting high values of m
 57 different fitness variables, $Y^{(1)}, \dots, Y^{(m)} \in \mathbb{R}$. For example, we are interested in generating diverse
 58 enzyme sequences that are both highly thermostable (*i.e.*, have high melting temperature) and exhibit
 59 high catalytic activity for a reaction of interest. Therefore, we let $Y^{(1)}$ denote melting temperature
 60 and let $Y^{(2)}$ denote catalytic activity, k_{cat}/K_M , which quantifies how efficiently an enzyme converts
 61 a particular substrate into a particular product. Toward this goal, we propose sampling from the
 62 sequence distribution, p^* , that solves the following optimization problem:

$$\begin{aligned} & \arg \max_{p \in \mathcal{P}} H(p) \\ & \text{subject to } \mathbb{E}_p[f_1(x)] \geq \tau_1, \\ & \quad \dots \\ & \quad \mathbb{E}_p[f_m(x)] \geq \tau_m, \\ & \quad \text{support}(p) \subseteq \bigcap_{j=1}^m \text{TRUSTREGION}_j \end{aligned} \tag{1}$$

63 where \mathcal{P} is the set of all distributions over \mathcal{X} , $H(p)$ denotes the entropy of the distribution p ,
 64 $\tau_1, \dots, \tau_m \in \mathbb{R}$ are desired lower bounds on the expected predicted fitnesses of sequences drawn
 65 from p^* , and $\text{TRUSTREGION}_j \subseteq \mathcal{X}$ denotes the region of sequence space on which we trust the
 66 predictions of the j -th predictive model, f_j (see examples in Section 2.1). Sequences sampled from
 67 the solution to this problem will simultaneously (i) have high average predicted fitnesses and (ii)
 68 belong to the region of sequence space on which we trust the predictions of all the models.

69 The distribution, p^* , that solves Eq. (1) has a likelihood of the following form:

$$p^*(x) \propto \begin{cases} \exp(\sum_{j=1}^m \lambda_j \cdot f_j(x)) & \text{if } x \in \bigcap_{j=1}^m \text{TRUSTREGION}_j \\ 0 & \text{otherwise} \end{cases} \tag{2}$$

70 where $\lambda_1, \dots, \lambda_m > 0$ are Lagrange multipliers corresponding to the first m constraints in Eq. (1).
 71 Our proposed approach is therefore to sample from p^* , where we treat $\lambda_1, \dots, \lambda_m$ as hyperparameters
 72 that navigate a trade-off between the expected predicted fitnesses of the designed sequences and their
 73 diversity, as quantified by the entropy of the distribution: higher values of $\lambda_1, \dots, \lambda_m$ lead to higher
 74 expected predicted fitnesses, but also sequence distributions with lower entropy.

75 2.1 Predictive models and trust regions for thermostability and catalytic activity

76 To design thermostable, active enzymes, we used the following two predictive models and corre-
 77 sponding trust regions. See the Appendix for more details.

78 **Predictive model of catalytic activity.** The fitness effect of adding mutations to a WT sequence
79 has been shown to be positively correlated with the likelihood of the mutant under natural protein
80 sequence density (NPSD) models (15; 11; 16), particularly for fitness functions related to the natural
81 function of the WT. We wanted to investigate whether such likelihoods could therefore be used to
82 design active enzymes without needing experimental measurements of catalytic activity. We set
83 $f_1(x) = f_{\text{ProGen}}(x)$ to be the average log-likelihood of a sequence, x , under five PROGEN2 models
84 of varying sizes: PROGEN2-base, PROGEN2-medium, PROGEN2-large, PROGEN2-xlarge, and
85 PROGEN2-BFD90.

86 To avoid ambiguity, we will refer to f_{ProGen} as the PROGEN2 *log-likelihood*, and will refer to the
87 (log-)likelihood given in Eq. (2) as the (*log-*)*likelihood under the target or designed distribution*.

88 **Trust region for catalytic activity.** Correlations between fitness and NPSD model likelihoods
89 have been demonstrated primarily on datasets comprising sequences that are mutants of a wild type
90 sequence (WT) known to be at least moderately fit (15; 11; 16). Based on this evidence, the region of
91 sequence space on which we trust the PROGEN2 log-likelihood of an enzyme sequence to be corre-
92 lated with its activity, and to which we therefore assign $\text{TRUSTREGION}_1 = \text{TRUSTREGION}_{\text{ProGen}}$,
93 is the set of sequences within M mutations of an active WT enzyme sequence¹. We applied our
94 method to four ACYP and four lysozyme WT sequences, ranging from 91 to 177 residues in length
95 (see Appendix for details). We used two different settings of M for $\text{TRUSTREGION}_{\text{ProGen}}$: we first
96 set it to the maximum number of mutations, M_{PROSS} , introduced by the baseline method PROSS
97 (9)—which we describe shortly—then also ran our method with $M = M_{\text{PROSS}} + 7$ for ACYPs and
98 $M = M_{\text{PROSS}} + 10$ for lysozymes. These two different settings will expose the effect of more and
99 less conservative trust regions on the PROGEN2 log-likelihood.

100 At a high level, PROSS first restricts candidate mutations at each site to amino acids observed at
101 that site in homologous sequences, then filters out any predicted by Rosetta to destabilize the WT.
102 It then returns the nine combinations of the remaining candidates that score best under the Rosetta
103 energy function (9). PROSS has been used successfully to design numerous thermostable enzymes
104 and other proteins (9; 10; 19) and represents a proven, competitive baseline. Note that PROSS begins
105 by constraining the set of candidate mutations, due the computational bottleneck of modeling variants
106 with many mutations with Rosetta. In contrast, because the cost of making predictions with our data-
107 driven models is trivial, we can tractably evaluate orders of magnitude more sequences than PROSS,
108 including mutants with more mutations away from the WT and mutants whose individual mutations
109 alone may not be beneficial. Ongoing experimental characterization will reveal whether this more
110 permissive exploration of sequence space enables our method to identify even more thermostable,
111 active variants than PROSS.

112 **Predictive model of thermostability.** We set $f_2(x) = f_{\text{melting}}(x)$ to be the predicted melting tem-
113 perature of a sequence x , according to a regression model fit to the “meltome” dataset, which contains
114 the experimentally measured melting temperatures of more than 34k unique protein sequences (14).
115 The regression model was an ensemble of ten two-layer neural networks, trained on ESM-1b (20)
116 embeddings of each sequence, where we took the average of the ensemble’s predictions.

117 **Trust region for thermostability.** We set $\text{TRUSTREGION}_2 = \text{TRUSTREGION}_{\text{melting}}$ to contain all
118 sequences classified as in-distribution with the meltome sequences, according to an out-of-distribution
119 detector we built following the k -nearest neighbors strategy in (21). Briefly, we took the outputs of
120 the final hidden layer in the thermostability predictive model for all the meltome sequences, then
121 classified a sequence as in-distribution if the Euclidean distance between its final-layer output and the
122 k -th nearest meltome sequence final-layer output was less than a threshold, c . The hyperparameters k
123 and c were set using three out-of-distribution sequence datasets (see the Appendix for details).

124 2.2 Metropolis-Hastings algorithm

125 We used a Metropolis-Hasting algorithm to sample from the distribution, p^* , with the
126 likelihood given in Eq. (2) (see the Appendix for details). We are working to-
127 ward open-sourcing code implementing our approach, which will soon be available at
128 <https://github.com/salesforce/designing-thermostable-enzymes>.

¹See (17; 18) for emerging evidence of this correlation on other kinds of datasets. However, to be conservative, our trust region mimics the “WT-centered” datasets on which this correlation has been much more robustly documented.

129 **3 *In silico* results**

130 **Predicted thermostability.** For all eight WT sequences, our method designed variants with higher
 131 predicted melting temperatures than either the WT or PROSS-designed variants (by at least 2°C to
 132 20°C over the latter), when constrained to the same number of mutations, M_{PROSS} , as the maximum
 133 number introduced by PROSS (Fig. 1, yellow and green diamonds compared to orange circles).
 134 When the constraint was relaxed, to $M_{\text{PROSS}} + 7$ for ACYPs and $M_{\text{PROSS}} + 10$ for lysozymes, our
 135 designed sequences achieved gains of at least 15°C to 30°C in predicted melting temperature over
 136 PROSS (Fig. 1, pink and purple diamonds compared to orange circles).

137 **PROGEN2 log-likelihood.** Examining the PROGEN2 log-likelihoods adds nuance to these results.
 138 Among the eight WTs, we observed two different *in silico* phenomena: for six WTs, both our method
 139 and PROSS designed variants with higher predicted melting temperatures and higher PROGEN2
 140 log-likelihoods (Fig. 1, left column), while for two WTs, our method designed variants predicted to
 141 be more thermostable but at the cost of lower PROGEN2 log-likelihood (Fig. 1, right column). In the
 142 latter cases, the PROSS-designed variants also had lower PROGEN2 log-likelihoods and/or lower
 143 predicted melting temperatures (Fig. 1, right column), which suggests some reassuring consistency
 144 between our approach and PROSS: when our method cannot improve predicted thermostability
 145 without sacrificing PROGEN2 log-likelihood, neither can PROSS. Note that the two challenging WTs,
 146 human ACYP-2 and T4 lysozyme (Fig. 1, right column), are known to be particularly active (22). It
 147 is therefore plausible that they are close to a Pareto frontier between activity and thermostability (8).

148 **Concluding remarks.** Wet-lab characterization of the sequences designed using both our approach
 149 and PROSS are currently underway. We look forward to sharing the results to advance the field’s
 150 understanding of how to employ data-driven models to design proteins enhanced for multiple fitness
 151 functions.

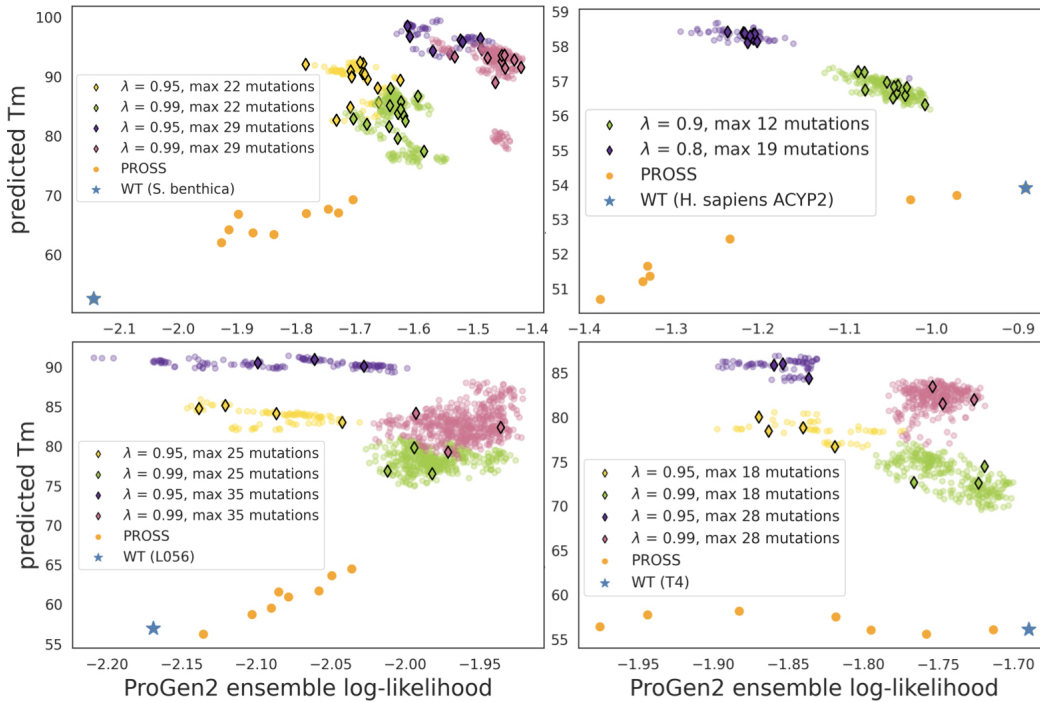


Figure 1: PROGEN2 log-likelihoods and predicted melting temperatures (C°) of variants designed using our approach and PROSS, for four example WT sequences. Top row: results for the ACYP WTs in *S. benthica* (left) and human (right). Bottom row: results for a designed lysozyme, L056 (12) (left), and the phage T4 lysozyme (right). In all subplots, the WT is the blue star, PROSS variants are orange circles, and the remaining circles are our designed sequences, where the colors correspond to different hyperparameter settings. The diamonds are random samples selected as our final designed variants, whose experimental characterizations are currently underway. The legend references two hyperparameters: the maximum number of mutations used for TRUSTREGION_{ProGen}, and λ , for which higher values correspond to higher PROGEN2 log-likelihoods.

References

- 152
- 153 [1] S. Biswas, G. Khimulya, E. C. Alley, K. M. Esvelt, and G. M. Church, “Low-N protein
154 engineering with data-efficient deep learning,” *Nature Methods*, vol. 18, no. 4, pp. 389–396,
155 2021.
- 156 [2] B. L. Hie and K. K. Yang, “Adaptive machine learning for protein engineering,” *Curr. Opin.*
157 *Struct. Biol.*, vol. 72, pp. 145–152, Dec. 2021.
- 158 [3] Z. Wu, S. B. J. Kan, R. D. Lewis, B. J. Wittmann, and F. H. Arnold, “Machine learning-assisted
159 directed protein evolution with combinatorial libraries,” *Proc. Natl. Acad. Sci. U. S. A.*, vol. 116,
160 pp. 8852–8858, Apr. 2019.
- 161 [4] K. K. Yang, Z. Wu, and F. H. Arnold, “Machine-learning-guided directed evolution for protein
162 engineering,” *Nat. Methods*, vol. 16, pp. 687–694, Aug. 2019.
- 163 [5] D. H. Brookes, H. Park, and J. Listgarten, “Conditioning by adaptive sampling for robust design,”
164 in *Proc. of the International Conference on Machine Learning (ICML)*, 2019.
- 165 [6] F. Rigoldi, S. Donini, A. Redaelli, E. Parisini, and A. Gautieri, “Review: Engineering of
166 thermostable enzymes for industrial applications,” *APL Bioeng*, vol. 2, p. 011501, Mar. 2018.
- 167 [7] M. M. Pinney, D. A. Mokhtari, E. Akiva, F. Yabukarski, D. M. Sanchez, R. Liang, T. Doukov,
168 T. J. Martinez, P. C. Babbitt, and D. Herschlag, “Parallel molecular mechanisms for enzyme
169 temperature adaptation,” *Science*, vol. 371, Mar. 2021.
- 170 [8] W. A. Baase, L. Liu, D. E. Tronrud, and B. W. Matthews, “Lessons from the lysozyme of phage
171 T4,” *Protein Sci.*, vol. 19, pp. 631–641, Apr. 2010.
- 172 [9] A. Goldenzweig, M. Goldsmith, S. E. Hill, O. Gertman, P. Laurino, Y. Ashani, O. Dym, T. Unger,
173 S. Albeck, J. Prilusky, R. L. Lieberman, A. Aharoni, I. Silman, J. L. Sussman, D. S. Tawfik,
174 and S. J. Fleishman, “Automated structure- and Sequence-Based design of proteins for high
175 bacterial expression and stability,” *Mol. Cell*, vol. 63, pp. 337–346, July 2016.
- 176 [10] Y. Peleg, R. Vincentelli, B. M. Collins, K.-E. Chen, E. K. Livingstone, S. Weeratunga, N. Leneva,
177 Q. Guo, K. Remans, K. Perez, G. E. K. Bjerga, Ø. Larsen, O. Vaněk, O. Skořepa, S. Jacquemin,
178 A. Poterszman, S. Kjær, E. Christodoulou, S. Albeck, O. Dym, E. Ainbinder, T. Unger,
179 A. Schuetz, S. Matthes, M. Bader, A. de Marco, P. Storici, M. S. Semrau, P. Stolt-Bergner,
180 C. Aigner, S. Suppmann, A. Goldenzweig, and S. J. Fleishman, “Community-Wide experimental
181 evaluation of the PROSS Stability-Design method,” *J. Mol. Biol.*, vol. 433, p. 166964, June
182 2021.
- 183 [11] E. Nijkamp, J. Ruffolo, E. N. Weinstein, N. Naik, and A. Madani, “ProGen2: Exploring the
184 boundaries of protein language models.” arXiv preprint 2206.13517, 2022.
- 185 [12] A. Madani, B. Krause, E. R. Greene, S. Subramanian, B. P. Mohr, J. M. Holton, J. L. Olmos,
186 C. Xiong, Z. Z. Sun, R. Socher, J. S. Fraser, and N. Naik, “Deep neural language modeling
187 enables functional protein generation across families.” bioRxiv preprint 2021.07.18.452833,
188 July 2021.
- 189 [13] B. L. Hie, D. Xu, V. R. Shanker, T. U. J. Bruun, P. A. Weidenbacher, S. Tang, and P. S. Kim,
190 “Efficient evolution of human antibodies from general protein language models and sequence
191 information alone.” bioRxiv preprint 2022.04.10.487811, 2022.
- 192 [14] A. Jarzab, N. Kurzawa, T. Hopf, M. Moerch, J. Zecha, N. Leijten, Y. Bian, E. Musiol,
193 M. Maschberger, G. Stoehr, I. Becher, C. Daly, P. Samaras, J. Mergner, B. Spanier, A. Angelov,
194 T. Werner, M. Bantscheff, M. Wilhelm, M. Klingenspor, S. Lemeer, W. Liebl, H. Hahne, M. M.
195 Savitski, and B. Kuster, “Meltome atlas—thermal proteome stability across the tree of life,” *Nat.*
196 *Methods*, Apr. 2020.
- 197 [15] J. Meier, R. Rao, R. Verkuil, J. Liu, T. Sercu, and A. Rives, “Language models enable zero-shot
198 prediction of the effects of mutations on protein function,” in *Advances in Neural Information*
199 *Processing Systems* (M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan,
200 eds.), vol. 34, pp. 29287–29303, Curran Associates, Inc., 2021.

- 201 [16] P. Notin, M. Dias, J. Frazer, J. Marchena-Hurtado, A. Gomez, D. S. Marks, and Y. Gal,
202 “Tranception: protein fitness prediction with autoregressive transformers and inference-time
203 retrieval,” in *Proc. of the International Conference on Machine Learning (ICML)*, 2022.
- 204 [17] W. P. Russ, M. Figliuzzi, C. Stocker, P. Barrat-Charlaix, M. Socolich, P. Kast, D. Hilvert,
205 R. Monasson, S. Cocco, M. Weigt, and R. Ranganathan, “An evolution-based model for
206 designing chorismate mutase enzymes,” *Science*, vol. 369, pp. 440–445, July 2020.
- 207 [18] C. Dallago, J. Mou, K. E. Johnston, B. J. Wittmann, N. Bhattacharya, S. Goldman, A. Madani,
208 and K. K. Yang, “Flip: Benchmark tasks in fitness landscape inference for proteins.” bioRxiv
209 preprint 2021.11.09.467890, 2021.
- 210 [19] A. Goldenzweig and S. J. Fleishman, “Principles of protein stability and their application in
211 computational design,” *Annu. Rev. Biochem.*, vol. 87, pp. 105–129, June 2018.
- 212 [20] A. Rives, J. Meier, T. Sercu, S. Goyal, Z. Lin, J. Liu, D. Guo, M. Ott, C. L. Zitnick, J. Ma, and
213 R. Fergus, “Biological structure and function emerge from scaling unsupervised learning to 250
214 million protein sequences,” *Proc. Natl. Acad. Sci. U. S. A.*, vol. 118, Apr. 2021.
- 215 [21] Y. Sun, Y. Ming, X. Zhu, and Y. Li, “Out-of-distribution detection with deep nearest neighbors,”
216 in *Proc. of the International Conference on Machine Learning (ICML)*, 2022.
- 217 [22] A. V. Gribenko, M. M. Patel, J. Liu, S. A. McCallum, C. Wang, and G. I. Makhatadze, “Rational
218 stabilization of enzymes by computational redesign of surface charge-charge interactions,” *Proc.*
219 *Natl. Acad. Sci. U. S. A.*, vol. 106, pp. 2601–2606, Feb. 2009.
- 220 [23] P. Leuenberger, S. Ganscha, A. Kahraman, V. Cappelletti, P. J. Boersema, C. von Mering,
221 M. Claassen, and P. Picotti, “Cell-wide analysis of protein thermal unfolding reveals determi-
222 nants of thermostability,” *Science*, vol. 355, Feb. 2017.
- 223 [24] E. J. Walker, J. Q. Bettinger, K. A. Welle, J. R. Hryhorenko, and S. Ghaemmaghami, “Global
224 analysis of methionine oxidation provides a census of folding stabilities for the human proteome,”
225 *Proc. Natl. Acad. Sci. U. S. A.*, vol. 116, pp. 6081–6090, Mar. 2019.

226 4 Appendix

227 4.1 Wild type sequences

228 We applied our method to design variants of the WT ACYP sequences in humans, the thermophile
229 *Pyrococcus horikoshii*, the psychrophile *Shewanella benthica*, and the ACYP-like domain in the hypF
230 WT in *Caldanaerobacter subterraneus*, as well as the WT lysozyme in phage T4 (8), two active
231 lysozymes denoted L056 and L070 designed using an NPSD model (12), and a glycoside hydrolase
232 WT found in *Brucella intermedia*. These WTs were chosen to span a wide range in both activity and
233 thermostability, and had lengths ranging from 91 to 108 residues for ACYPs and 164 to 177 residues
234 for lysozymes.

235 4.2 Predictive model of thermostability

236 The melting temperature predictive model can provide reasonably accurate predictions on a held-out
237 enzyme family (Fig. 2). Forthcoming experimental data will clarify whether this predictive accuracy
238 is sufficient to design enzyme variants with increased melting temperatures.

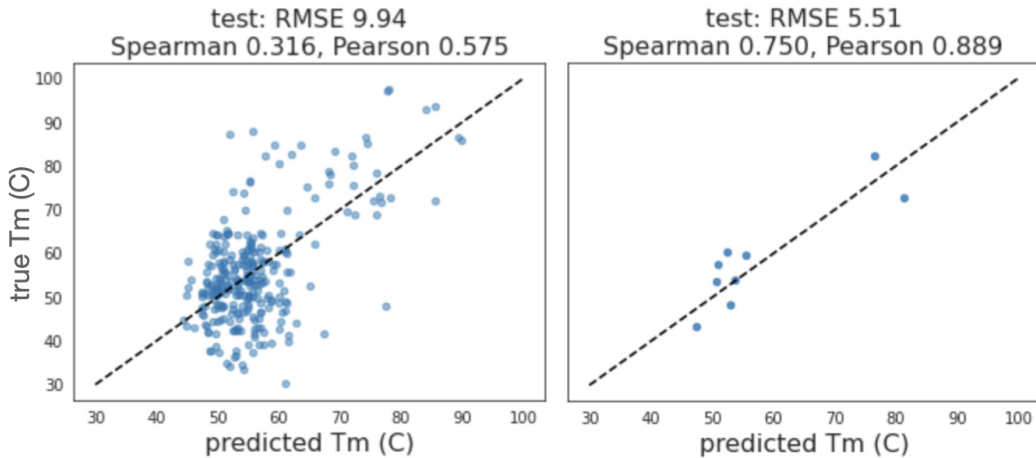


Figure 2: Performance of the thermostability predictive model on held-out ACYP-like sequences. Left: To assess the model’s performance when the training and test sequences are very different, the meltome sequences were partitioned into two sets: those with sequence identity of less than 50% with each of the nine ACYP sequences in the meltome, and those with sequence identity of 50% or more with any of these ACYPs. The model was trained on the former, and predictions on the latter are shown. Right: The model was trained on all meltome sequences other than the nine ACYP sequences, and predictions on those nine ACYPs are shown. This was the final predictive model used to design ACYP variants.

239 4.3 Trust region for thermostability

240 To build the out-of-distribution (OOD) detector, following (21), we computed the final hidden layer
241 outputs of the thermostability predictive model for all sequences in the training data—that is, all 34k
242 unique protein sequences in the meltome dataset, other than the nine excluded ACYP sequences. Let
243 $z_i = \phi(x_i) / \|\phi(x_i)\|_2$ denote the normalized final-layer representation of the i -th training sequence,
244 x_i . Then given a test sequence, x' , let $z_{(i)}$ denote the training sequence representation that is i -th
245 closest to $z' = \phi(x') / \|\phi(x')\|_2$ in Euclidean distance. The test sequence is then classified as OOD if
246 $\|z' - z_{(k)}\|_2 > c$, where k and c are hyperparameters that we set as follows. Following (21), for any
247 value of k , the threshold, c , was set to the smallest value such that 95% of the training sequences are
248 classified as in-distribution. The integer k was then set to the value that maximized the rate of OOD
249 classification on sequences from three different OOD sequence datasets: (i) domain-length peptide
250 sequences from (23), (ii) domain-length peptide sequences from (24), and (iii) full-length enzyme
251 sequences from (7). The resulting hyperparameter values were $k = 10$ and $c = 3 \times 10^{-5}$.

252 **4.4 Hyperparameter settings**

253 The hyperparameters $\lambda_1 = \lambda_{\text{Progen}}$ and $\lambda_2 = \lambda_{\text{melting}}$ in Eq. (2) were set as $\lambda_{\text{Progen}} = \lambda_{\text{scale}}\lambda$ and
254 $\lambda_{\text{melting}} = \lambda_{\text{scale}}(1 - \lambda)$, where $\lambda_{\text{scale}} = 300$ and we used two different settings of λ for each WT
255 sequence. For all WTs except for human ACYP-2, we ran our method with $\lambda \in \{0.99, 0.95\}$. For
256 human ACYP-2, we used $\lambda \in \{0.8, 0.9\}$.

257 **4.5 Metropolis-Hastings algorithm**

258 We now describe the proposal distribution, $q(x' | x)$, that we developed for the Metropolis-Hastings
259 (MH) algorithm, where $q(x' | x)$ denotes the probability of proposing the sequence x' given the
260 current sequence, x , at each iteration of the MH algorithm. We found this proposal distribution to
261 achieve an effective trade-off between producing sequences with high likelihood under the target
262 distribution, thereby allowing the sampling algorithm to mix and converge in fewer iterations, and
263 being computationally lightweight. To motivate its construction, we first observed that with naive
264 proposal distributions, such as sampling mutants of the current sequence uniformly at random, it
265 was often more difficult to propose a sequence, x' , with higher $f_{\text{ProGen}}(x')$ than one with higher
266 $f_{\text{melting}}(x')$. Therefore, we developed a proposal distribution intended to more efficiently propose
267 mutants of the current sequence with higher $f_{\text{ProGen}}(x')$. The high-level intuition is that, given a
268 multiple sequence alignment (MSA) of homologous sequences of a WT sequence, the position weight
269 matrix (PWM) containing the frequencies of the twenty amino acids observed at each site should
270 provide some first-order information as to which nearby mutants have high PROGEN2 likelihoods.
271 However, we did not want practitioners to have to manually curate MSAs in order to use our method.
272 As an alternative, we found that for many WT sequences, the entries of the following *likelihood-based*
273 *position weight matrix* (L-PWM), $A \in \mathbb{R}^{L \times |A|}$, are strongly correlated with the PWM computed
274 from an MSA:

$$A_{la} = \frac{1}{Z_l} \exp(\gamma \cdot (f_{\text{ProGen}}(x_{l,a}) - f_{\text{ProGen}}(x_{\text{WT}}))),$$
$$Z_l = \sum_{a' \in \mathcal{A}} \exp(\gamma \cdot (f_{\text{ProGen}}(x_{l,a'}) - f_{\text{ProGen}}(x_{\text{WT}})))$$

275 where we set the hyperparameter $\gamma = 80$ to the value that maximized the average Pearson correlation
276 between the entries, A_{la} , and the corresponding entries of the MSA-based PWM across WTs from
277 several different protein families.

278 At each iteration of the MH algorithm, given the current sequence, x , we generate a proposal, x' ,
279 using the following procedure:

- 280 1. Each position in the sequence x has an independent probability, $p_{\text{mutate}} = 0.02$, of being
281 mutated.
- 282 2. For each position, $l \in \{1, \dots, L\}$, that is to be mutated, the probability of mutating to the
283 a -th amino acid is given by A_{la} .

284 We then accept or reject the proposal as prescribed by the MH algorithm. We ran the algorithm for 10k
285 iterations and observed that chains appeared to mix and converge by at most 6k iterations. Therefore,
286 we took a random sample of the sequences after 6k iterations as our final designed sequences.