

# Designing active and thermostable enzymes with sequence-only predictive models

Clara Fannjiang<sup>1,2</sup>, Micah Olivas<sup>3</sup>, Eric R. Greene<sup>4</sup>, Craig J. Markin<sup>5</sup>, Bram Wallace<sup>2</sup>, Ben Krause<sup>2</sup>, Margaux M. Pinney<sup>6</sup>, James S. Fraser<sup>4</sup>, Polly M. Fordyce<sup>3,7</sup>, Ali Madani<sup>8</sup>, Nikhil Naik<sup>2</sup>

<sup>1</sup>Dept. of Electrical Engineering & Computer Sciences, UC Berkeley <sup>2</sup>Salesforce Research <sup>3</sup>Dept. of Genetics, Stanford University  
<sup>4</sup>Dept. of Bioengineering & Therapeutic Sciences, UCSF <sup>5</sup>Dept. of Biochemistry, Stanford University <sup>6</sup>Dept. of Biochemistry & Biophysics, UCSF  
<sup>7</sup>Dept. of Bioengineering, Stanford University <sup>8</sup>Profluent Bio

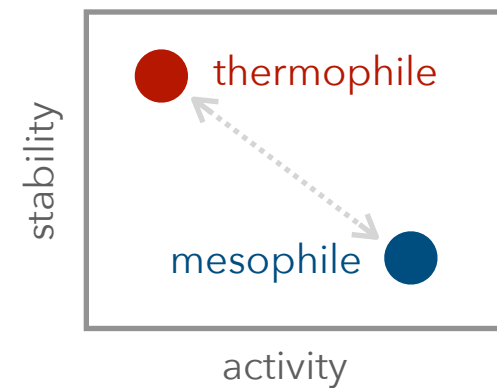


How can we use predictive models of fitness to design proteins

- (i) that satisfy *multiple* properties
- (ii) when these models are not always trustworthy?

## Case study: designing active, more thermostable enzymes

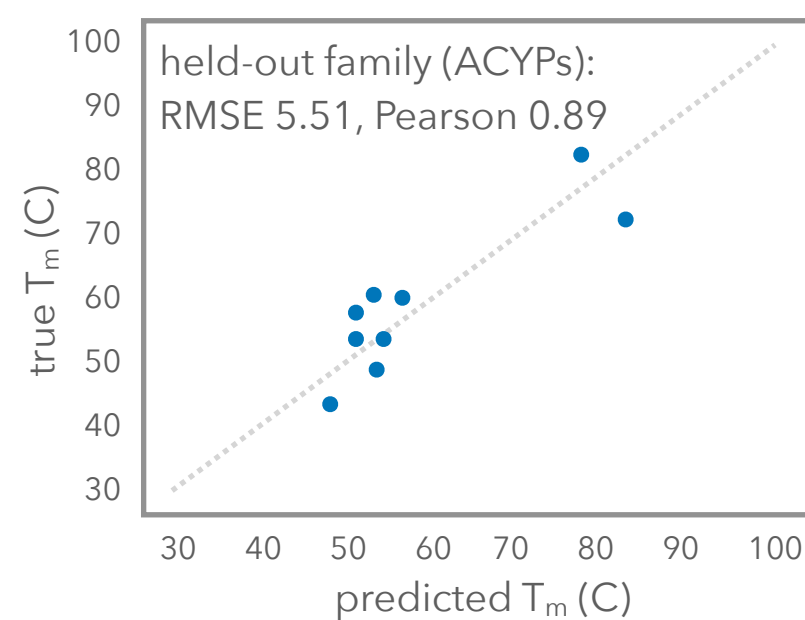
- broadly applicable goal, e.g. for industrial applications
- natural enzymes often exhibit trade-off
- existing methods: PROSS, consensus



## Predictive models and trust regions for...

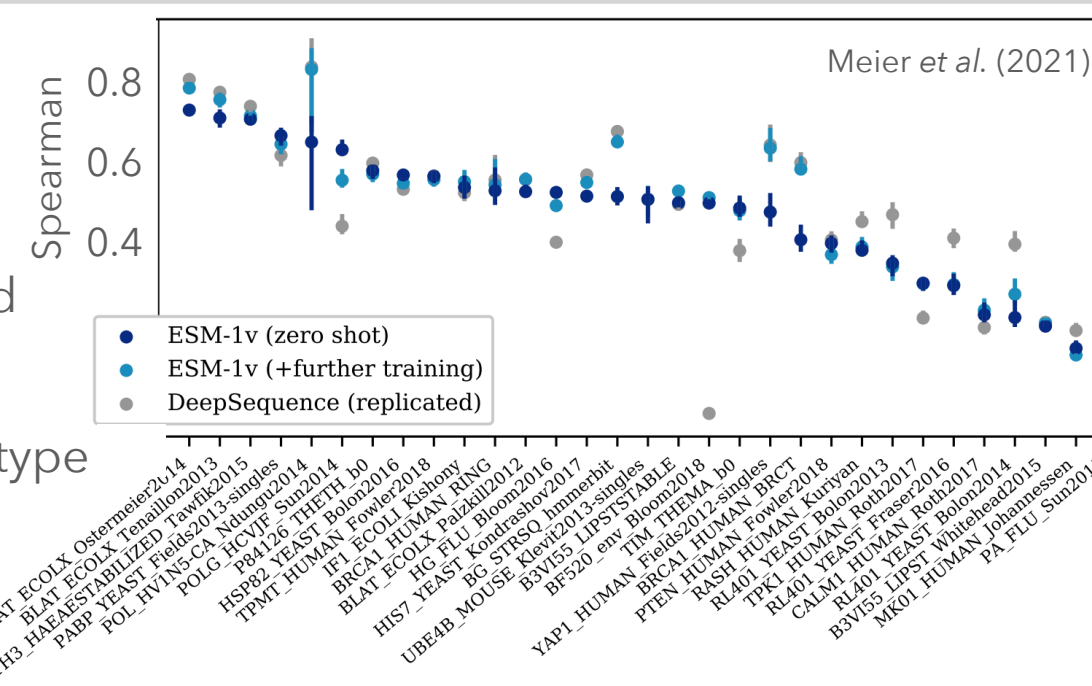
### Thermostability

- meltome: melting points for 34k+ protein sequences
- NN on top of ESM-1b embeddings
- trust region: sequences classified as in-distribution with meltome sequences



### Activity

- ProGen2 log-likelihood
- trust region: within  $M$  mutations from a wild type



We also expect the ProGen2 log-likelihood's correlation with activity to depend on the wild type—specifically, **on the extent to which activity drove evolutionary pressure on the wild type.**

## Our general approach

$f_i, i = 1, \dots, m$ : predictive model of  $i$ -th fitness function  
 $\text{TRUSTREGION}_i, i = 1, \dots, m$ : region of sequence space on which we trust  $f_i$

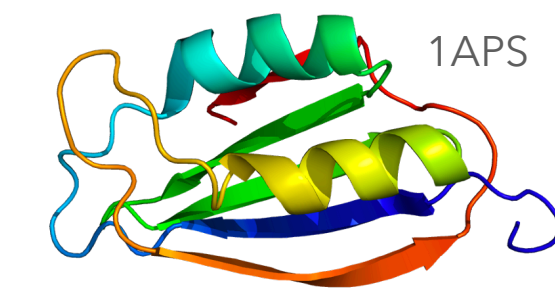
We use a Metropolis-Hastings algorithm to sample novel sequences from the distribution

$$p^*(x) \propto \begin{cases} \exp(\sum_{i=1}^m \lambda_i \cdot f_i(x)) & \text{if } x \in \bigcap_{i=1}^m \text{TRUSTREGION}_i \\ 0 & \text{otherwise} \end{cases}$$

which is the solution to the following optimization problem:

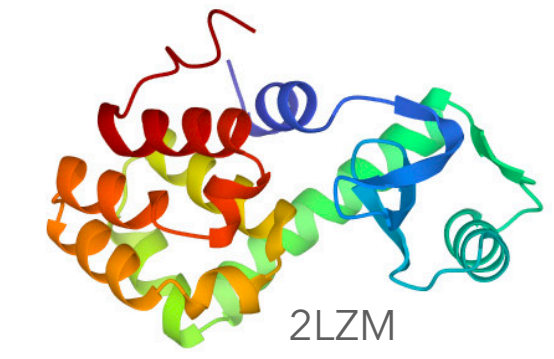
$$\begin{aligned} & \arg \max_{p \in \mathcal{P}} H(p) \\ & \text{subject to } \mathbb{E}_p[f_1(x)] \geq \tau_1, \\ & \quad \dots \\ & \mathbb{E}_p[f_m(x)] \geq \tau_m, \\ & \text{support}(p) \subseteq \bigcap_{j=1}^m \text{TRUSTREGION}_j \end{aligned}$$

## Wild-type enzymes



### Acylphosphatase (ACYP)

- human ACYP2
- *P. horikoshii* (thermophile)
- *S. benthica* (psychrophile)
- ACYP-like domain in hypF



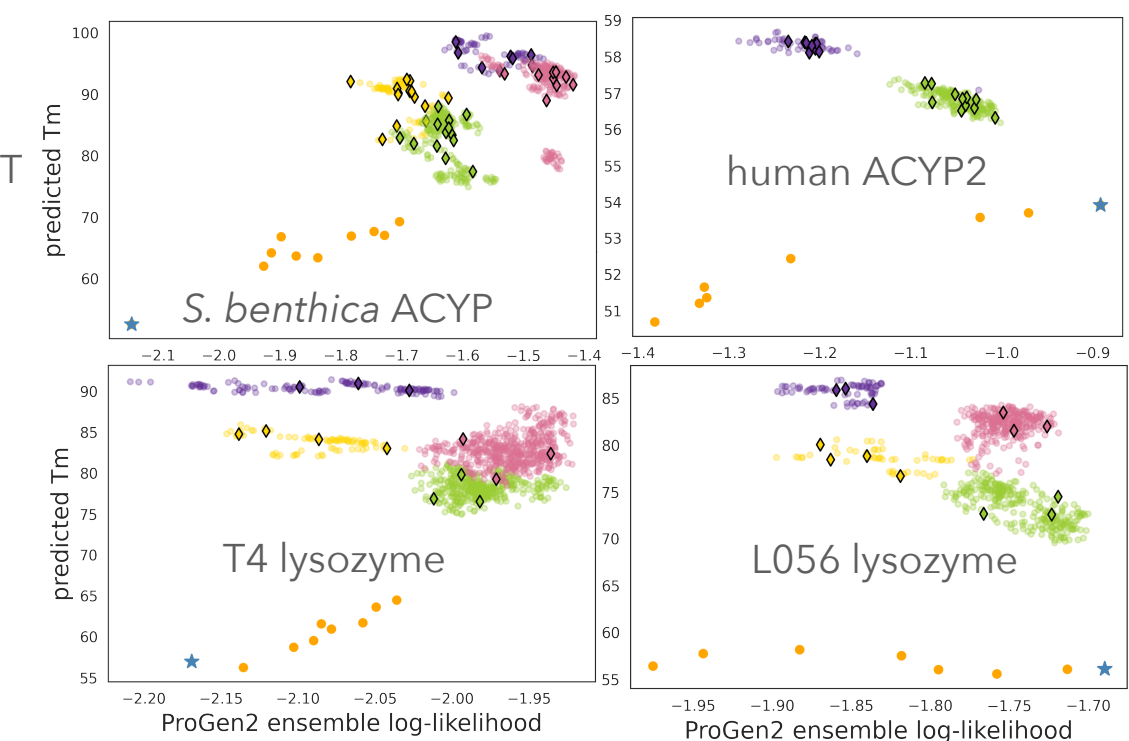
### Lysozyme

- phage T4
- L056 (previously designed)
- L070 (previously designed)
- *B. intermedia*

## Predicted behavior of designed sequences

Experimental characterization is in progress and forthcoming

- ♦ low  $\lambda$ , high ID to WT
- ♦ high  $\lambda$ , high ID
- ♦ low  $\lambda$ , low ID
- ♦ high  $\lambda$ , low ID
- PROSS
- ★ wild type (WT)



### References

Goldenzweig et al. (2016), *Mol. Cell*  
 Jarzab et al. (2019), *Nat. Methods*  
 Madani et al. (2022), *Nat. Biotech.* (to appear)  
 Markin & Mokhtari et al. (2021), *Science*  
 Meier et al. (2021), *NeurIPS*

Nijkamp & Ruffolo et al. (2022), arXiv:2206.13517  
 Pinney et al. (2021), *Science*  
 Rives et al. (2021), *PNAS*  
 Sun et al. (2022), *ICML*