

ProSpero: Active Learning for Robust Protein Design Beyond Wild-Type Neighborhoods

Introduction & Motivation

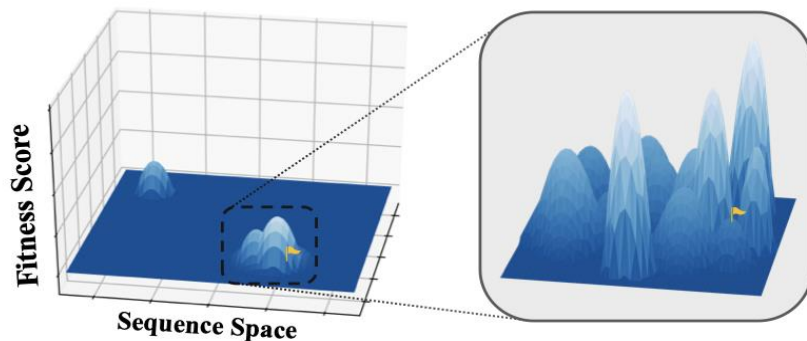


Introduction

Main goal – Designing novel protein sequences with desired properties

Challenges:

- rugged and sparse "fitness landscape"
- combinatorial search space
- expensive black-box evaluations



How to reliably explore further away from the wild-type?

Wild-type – A naturally occurring sequence serving as a reference/starting point in the optimization

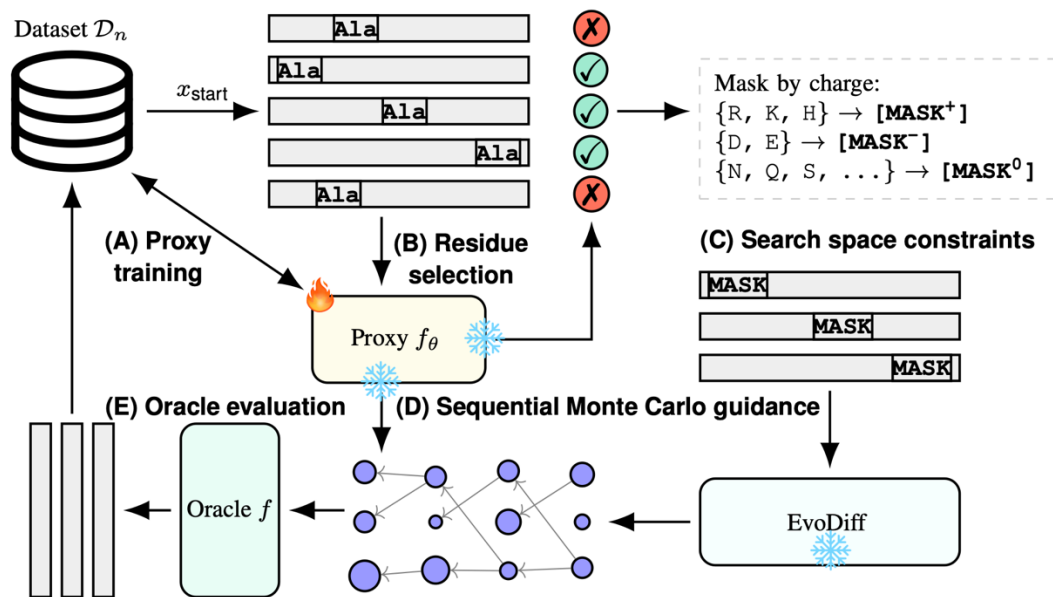
Potential solutions:

- Active learning: iterative re-training of the surrogate to progressively expand its support
- Biological priors: leverage prior biological knowledge to ensure plausibility even with a potentially misspecified surrogate

Proposed Framework

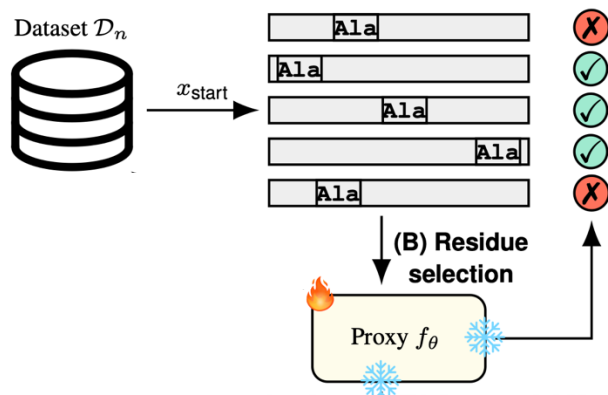


ProSpero



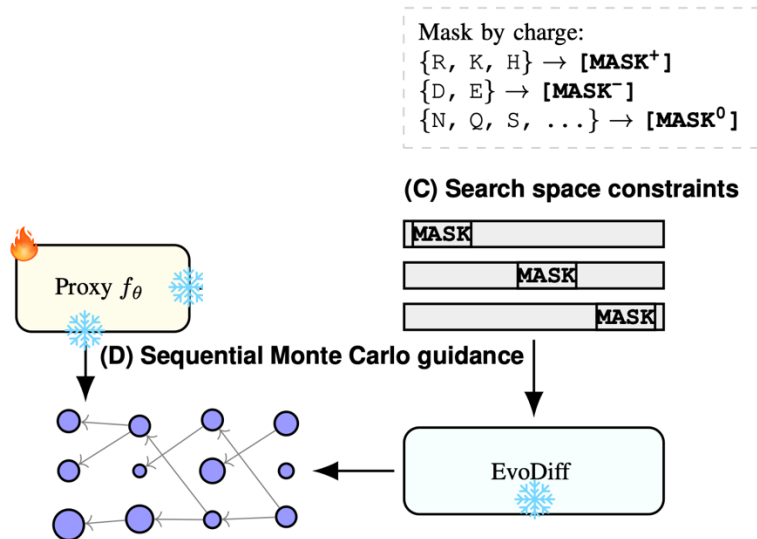
Inference-time guidance of a pre-trained pLM with a surrogate updated in an active learning loop – seamless integration of biological priors into online optimization, regardless of the target protein family

Targeted Masking



We focus edits on fitness-relevant residues, while preserving structurally and functionally important sites

Biologically-constrained Sequential Monte Carlo



Biologically-constrained Sequential Monte Carlo

Constrained proposal

$$\tilde{x}_{\pi(t)}^{(i)} \sim \mathcal{P}_{RAA}(\tilde{x}_{\pi(t)}^{(i)} \mid \tilde{x}_{\pi(<t)}^{(i)})$$

Weighting

$$x_{\text{unroll}}^{(i)} \sim \prod_{s=t+1}^T \mathcal{P}_{RAA}(\tilde{x}_{\pi(s)}^{(i)} \mid \tilde{x}_{\pi(<s)}^{(i)})$$

$$\hat{y}^{(i)} = \mu_{\theta}(x_{\text{unroll}}^{(i)}) + k \cdot \sigma_{\theta}(x_{\text{unroll}}^{(i)})$$

$$\text{Perp}(x_{\text{unroll}}^{(i)}) = \exp \left(-\frac{1}{|I|} \sum_{s=T-|I|+1}^T \log \mathcal{P}(x_{\text{unroll}_{\pi(s)}}^{(i)} \mid x_{\text{unroll}_{\pi(<s)}}^{(i)}) \right)$$

$$w_t^{(i)} = \hat{y}^{(i)} / \text{Perp}(x_{\text{unroll}}^{(i)})$$

Resampling

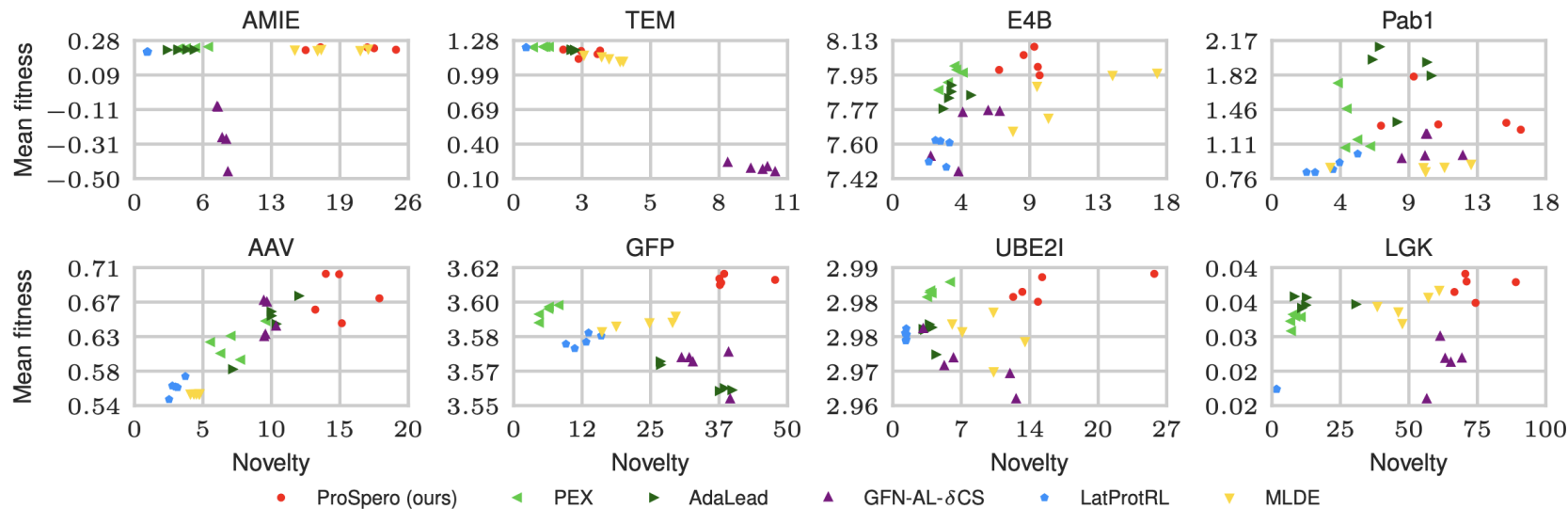
$$\tilde{x}^{(i)} \sim \text{Cat} \left(\{\tilde{x}^{(i)}\}_{i=1}^B, \left\{ \frac{w_t^{(i)}}{\sum_{j=1}^B w_t^{(j)}} \right\}_{i=1}^B \right)$$

Biologically-constrained SMC restricts proposals to residues with similar phys-chem properties to their wild-type counterparts, improving the likelihood of finding high fitness sequences under surrogate misspecification

Results



Breaking the fitness–novelty Pareto front



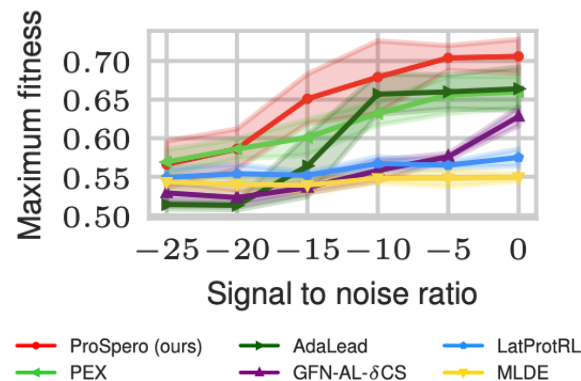
Robustness to surrogate misspecification

| Method | Maximum pTM | Mean pTM | Diversity | Novelty |
|---------------------|-------------------------------------|-------------------------------------|------------------------------------|------------------------------------|
| AdaLead | 0.796 ± 0.013 | 0.755 ± 0.011 | 8.83 ± 2.54 | 8.36 ± 2.97 |
| PEX | 0.807 ± 0.023 | 0.760 ± 0.012 | 6.14 ± 0.89 | 4.45 ± 0.38 |
| GFN-AL- δ CS | 0.791 ± 0.010 | 0.729 ± 0.005 | 16.92 ± 0.88 | 9.56 ± 0.60 |
| MLDE | 0.810 ± 0.020 | 0.752 ± 0.004 | 9.89 ± 1.11 | 20.88 ± 2.98 |
| LatProtRL | 0.787 ± 0.013 | 0.743 ± 0.003 | 6.32 ± 0.32 | 5.90 ± 0.53 |
| ProSpero | 0.822 ± 0.027 | 0.777 ± 0.020 | 11.50 ± 1.62 | 17.74 ± 3.20 |

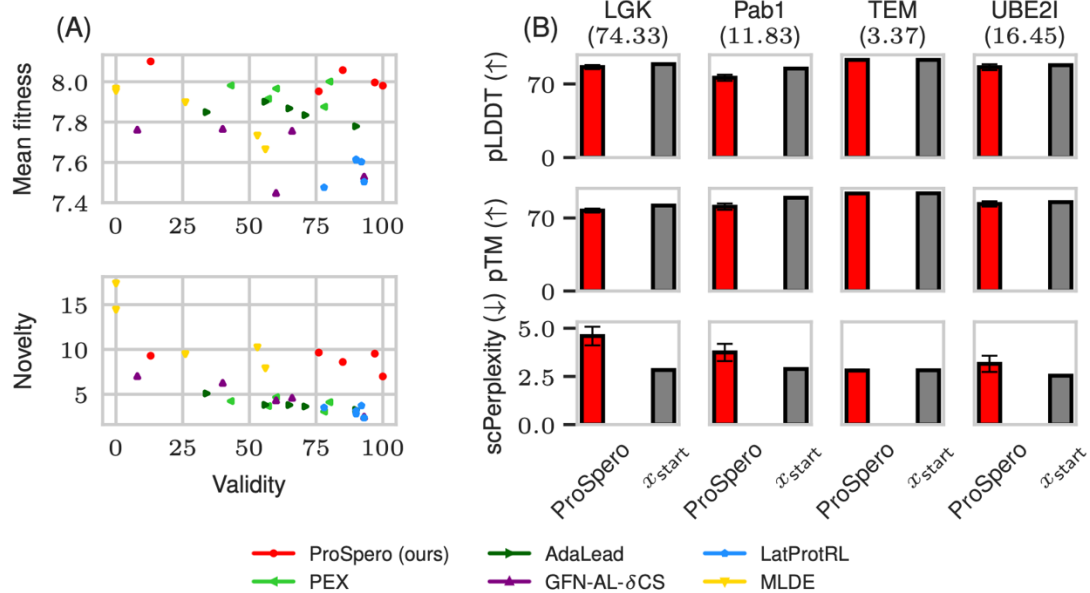
Starting sequence different from the wild-type by 35 residues

| Method | Maximum pTM | Mean pTM | Diversity | Novelty |
|---------------------|-------------------------------------|-------------------------------------|------------------------------------|------------------------------------|
| AdaLead | 0.593 ± 0.028 | 0.526 ± 0.007 | 14.26 ± 1.91 | 7.66 ± 1.08 |
| PEX | 0.578 ± 0.014 | 0.518 ± 0.003 | 3.40 ± 0.07 | 1.72 ± 0.04 |
| GFN-AL- δ CS | 0.630 ± 0.024 | 0.542 ± 0.006 | 24.13 ± 1.47 | 14.63 ± 1.16 |
| MLDE | 0.652 ± 0.059 | 0.572 ± 0.035 | 13.10 ± 1.18 | 21.68 ± 3.85 |
| LatProtRL | 0.560 ± 0.000 | 0.508 ± 0.003 | 2.24 ± 0.14 | 1.78 ± 0.16 |
| ProSpero | 0.672 ± 0.031 | 0.599 ± 0.014 | 14.51 ± 1.99 | 22.03 ± 1.69 |

Starting sequence different from the wild-type by 75 residues



Biologically plausible sequences



Conclusion



Conclusion

ProSpero facilitates protein design beyond wild-type neighborhoods by incorporating biological priors through:

- Inference-time guidance of a pre-trained pLM with a surrogate updated in an active learning loop
- Targeted masking of fitness-relevant residues while preserving key structural sites
- Biologically-constrained SMC sampling that restricts proposals to wild-type-like residues

We demonstrated robustness of ProSpero across diverse *in silico* protein engineering tasks



Thank you!

Ewa Szczurek



Vincent Fortuin

