

Multi-Objective Active Learning for Nanobody Development

Katharina Dost *1, <u>Klara Kropivšek*2</u>, Christian L. Camacho Villalón¹, Sašo Džeroski¹, and Ario de Marco²

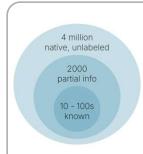
Motivation

Nanobodies:

- small "keys" that match specific sites on cells, proteins
- can bind to block or modulate specific functions
- can carry payload (e.g., drug, fluorescent marker)

Advantages:

- Small size → reach hidden epitopes that classical antibodies cannot reach
- Stable & robust → more tolerant to harsh conditions
- Versatile → can be engineered for therapy, diagnostics, research



Dataset

Small dataset (10 - 100s): in-house nanobodies with measured yields & developability (high-quality labels).

Medium dataset (~2,000): nanobodies with partial annotations (literature/structural).

Largest dataset (4 million): nonredundant native repertoire sequences (no experimental data, unlabeled).

Research Question

Obtaining experimental data is **costly**. A great ML model could guide discovery as it can point to potential candidates for a task, ensuring **developability** in the wet lab. A great model needs **data** to learn from.

Which experiments (= labeled data) would help to train such a model?

Big-Picture Idea

Initial Dataset (embedded with Abland)

Add labeled Find most informative nanobodies to dataset Obtain labels (lab experiments)

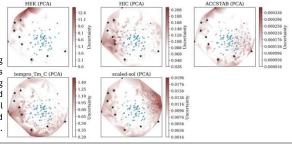
Active Learning

Uncertainty

Test Set

Labeled Dataset

Uncertainty Map highlighting areas with uncertain predictions per target in Ablang embedding space (dimensionality reduced with PCA). Blue = initial nanobodies; black = selected batch of nanobodies.



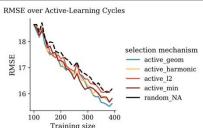
Active Learning Cycle

In each AL cycle, to select batch of nanobodies (e.g., 10):

- filter based on **constraints** for yield, thermal stability, solubility (as predicted)
- drop outliers
- select top X nanobodies that maximize uncertainty in predicting
- yield, polyreactivity, accelerated & thermal stability, solubility
- greedily select diverse subset

Simulation

- Simulation on 3000 nanobodies from our native library
- Labels were obtained using Abpred², TEMPRO³ and Protein-Sol⁴
- Per-target models with pertarget uncertainties; different aggregation strategies



Conclusion

Summary:

 Multi-objective active nanobody selection strategy for better property predictions

Next Steps:

- Model catering to distribution shift
- Refined multi-objective optimization techniques
- Tests on complete native library
- Real AL cycle + tags

Acknowledgements. Dost, Kropivšek and Camacho are supported by the European Union's Horizon Europe research and innovation programme under the Marie Skłodowska-Curie Postdoctoral Fellowship Programme, SMASH, co-funded under the grant agreement No. 101081355. The SMASH project is co-funded by the Republic of Slovenia and the European Union from the European Regional Development Fund. Džeroski and de Marco are supported by the Slovenian Research and Innovation Agency (under grants P2-0103, GC-0001, and P3-0428, N4-0282, N4-0325, J4-50144, respectively).

References

- ¹ https://github.com/oxpig/AbLang2
- ² https://github.com/maxhebditch/abpred
- https://github.com/Jerome-Alvarez/TEMPRO
 https://protein-sol.manchester.ac.uk/







Scan me to



¹ Jožef Stefan Institute, Slovenia

² Laboratory for Environment and Life Sciences, University of Nova Gorica, Slovenia